



Developmental Gene Discovery in a Hemimetabolous Insect: De Novo Assembly and Annotation of a Transcriptome for the Cricket *Gryllus Bimaculatus*

Citation

Zeng, Victor, Benjamin Scott Ewen-Campen, Hadley W. Horsch, Siegfried Roth, Taro Mito, and Cassandra G. Extavour. 2013. Developmental gene discovery in a hemimetabolous insect: De novo assembly and annotation of a transcriptome for the cricket *Gryllus bimaculatus*. PLoS ONE 8(5): e61479.

Published Version

doi:10.1371/journal.pone.0061479

Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:10498892>

Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP>

Share Your Story

The Harvard community has made this article openly available.
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

Developmental gene discovery in a hemimetabolous insect: *de novo* assembly and annotation of a transcriptome for the cricket *Gryllus bimaculatus*

Victor Zeng^{1,2}, Ben Ewen-Campen¹, Hadley W. Horch³, Siegfried Roth⁴, Taro Mito⁵, Cassandra G. Extavour^{1, 6}

1. Department of Organismic and Evolutionary Biology, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA.

2. Current address: Stylux Incorporated, 25 Stickney Road, Atkinson, NH 03811, USA

3. Departments of Biology and Neuroscience, Bowdoin College, 130B Druckenmiller Hall, Brunswick, ME 04011, USA.

4. Institute for Developmental Biology, University of Cologne, Cologne Biocenter, Zùlpicher StraÙe 47b, 50674, Cologne, Germany.

5. Department of Life Systems, Institute of Technology and Science, The University of Tokushima Graduate School, 2-1 Minami-Jyosanjima-cho, Tokushima City, 770-8506, Japan.

6. Author for correspondence. Email extavour@oeb.harvard.edu; Tel. (617) 496 1935; Fax (617) 496-9507

Author Email Addresses:

VZ victor.zeng@styluxdesigns.com

BEC bewencampen@oeb.harvard.edu

HH hhorch@bowdoin.edu

SR Siegfried.Roth@uni-koeln.de

TM mito@bio.tokushima-u.ac.jp

CGE extavour@oeb.harvard.edu

Abstract

Most genomic resources available for insects represent the Holometabola, which are insects that undergo complete metamorphosis like beetles and flies. In contrast, the Hemimetabola (direct developing insects), representing the basal branches of the insect tree, have very few genomic resources. We have therefore created a large and publicly available transcriptome for the hemimetabolous insect *Gryllus bimaculatus* (cricket), a well-developed laboratory model organism whose potential for functional genetic experiments is currently limited by the absence of genomic resources. cDNA was prepared using mRNA obtained from adult ovaries containing all stages of oogenesis, and from embryos samples on each day of embryogenesis. Using 454 Titanium pyrosequencing, we sequenced over four million raw reads, and assembled them into 21,512 isotigs (predicted transcripts) and 120,805 singletons with an average coverage per base pair of 51.3. We annotated the transcriptome manually for over 400 conserved genes involved in embryonic patterning, gametogenesis, and signaling pathways. BLAST comparison of the transcriptome against the NCBI non-redundant protein database (**nr**) identified significant similarity to **nr** sequences for 55.5% of transcriptome sequences, and suggested that the transcriptome may contain 19,874 unique transcripts. For predicted transcripts without significant similarity to known sequences, we assessed their similarity to other orthopteran sequences, and determined that these transcripts contain recognizable protein domains, largely of unknown function. We created a searchable, web-based database to allow public access to all raw, assembled and annotated data. This database is to our knowledge the largest *de novo* assembled and annotated transcriptome

resource available for any hemimetabolous insect. We therefore anticipate that these data will contribute significantly to more effective and higher-throughput deployment of molecular analysis tools in *Gryllus*.

Keywords: Hemimetabola; Orthoptera; cricket; 454 pyrosequencing; *de novo* transcriptome; Domain of Unknown Function (DUF)

Introduction

The vast majority of existing insect genomic resources are for the Holometabola or “higher insects,” which undergo true metamorphosis. These include disease vectors such as the mosquito *Anopheles gambiae* [1], agricultural pests such as the flour beetle *Tribolium castaneum* [2], and the powerful genetic model organism *Drosophila melanogaster* [3,4]. However, there are very few complete genome sequences available for the Hemimetabola or “lower insects”, which do not undergo true metamorphosis and branch basally to the Holometabola. Only three of the over 146,000 estimated species of hemimetabolous insects [5] have available genome sequences: the aphid *Acyrtosiphon pisum* [6], the kissing bug *Rhodnius prolixus* [7,8], and the human body louse *Pediculus humanus* [9]. Moreover, sequence divergence is so great among insects [10] that a specific genome cannot be used as a reference sequence for other insects even within the same order [see for example 11].

Among the Hemimetabola, the basally branching orthopteroid orders of insects are of particular interest to many fields of biology. Orthopterans have served as classical model organisms for neurobiology for several decades [12]. Multiple cricket species have been used for important studies of ecologically relevant polyphenisms [reviewed in 13], the evolution of endocrine functions and photobiology [14,15,16,17], speciation [18,19,20,21,22] and the evolution of behavior [23,24,25]. Crickets and locusts have also been important for addressing outstanding questions in evolutionary developmental biology, such as the evolution of molecular mechanisms for regeneration, segmentation, and axial patterning [26,27,28,29,30,31,32,33]. However, *de novo* genome assembly for organisms with extremely large genome sizes is costly and challenging [34,35,36].

Grasshopper genomes can be over twice as large as the human genome [37], and even the genome of the laboratory model cricket *Gryllus bimaculatus* is estimated at 1.7 Gbp (C. G. Extavour and R. Gregory, unpublished). If orthopteran genome projects are eventually undertaken, their annotation success will be significantly enhanced by the availability of large transcriptomes, but these are also few in number.

To date, only three Sanger-based EST projects and one large *de novo* assembled transcriptome generated with next-generation sequencing have been reported for orthopterans (Table 1). These projects have focused on specific post-embryonic developmental stages of pest locusts (*L. migratoria*, *S. gregaria*) and on the CNS of a cricket (*L. kohalensis*). Although most functional genetic studies on orthopterans focus on embryonic development [e.g. 28,29,38,39] and neurophysiological studies are increasingly examining the embryonic origins of neural structures and functions [e.g. 16,40,41,42,43], a transcriptome enriched for embryonic developmental transcripts is lacking. Here we present such a transcriptome for the model laboratory cricket, *G. bimaculatus*.

G. bimaculatus is a highly tractable orthopteran model for functional genetic studies in the laboratory (Fig. 1). Gene knockdown can be achieved by RNA interference during embryonic, post-embryonic and regenerative development [32,43,44]. *G. bimaculatus* is also the only orthopteran for which stable germ line transgenesis has been established [39]. Moreover, protocols for targeted genome editing using zinc finger nucleases or TALE nucleases have recently been developed [45]. However, all *G. bimaculatus* genes studied to date have been obtained by degenerate PCR [e.g. 28,46] or from limited Sanger-based EST libraries that are not available in an annotated database [e.g. 26].

In this report we present a *de novo* assembled and annotated transcriptome for *G. bimaculatus* oogenesis and embryonic development. We show that this transcriptome contains more putative unique gene transcripts than previous orthopteran transcriptomes, and adds sequence data to known GenBank accessions for *G. bimaculatus*. We manually annotate over 400 developmental genes, and develop an automated annotation method for the entire transcriptome based on similarity to *Drosophila* sequences. For predicted transcripts that lack significant similarity to GenBank accessions, we examine specifically those that are more similar to known orthopteran sequences, and find that the most represented predicted protein domains of such “orthopteroid” transcripts are domains of unknown function (DUFs). In contrast, the most represented predicted protein domains of transcripts of the transcriptome overall are zinc finger domains. Finally, we created a publicly accessible repository and database for the transcriptome, which is searchable by BLAST, pre-computed BLAST hits, or putative orthology assignments (gene names) derived from both manual and automated annotation.

Materials and Methods

Animal culture and collection of tissues for cDNA synthesis

G. bimaculatus cultures were maintained as previously described [28], at 28-29°C on a diet of oatmeal, wheat germ, soya protein, corn meal, sugar, yeast, salt, corn oil and Purina Cat Chow. This non-isogenic culture derives from a population of *G. bimaculatus* obtained from Livefoods Direct (Sheffield, UK), and was maintained as an inbred, self-sustaining culture for four years (or approximately 26 generations) prior to tissue

collection. We do not have estimates of genetic polymorphism for this population, so that accurate interpretation of putative SNP data is not possible in the present analysis.

Separate egg collections (total mass 781 mg) of 50-100 embryos on each of the first eight days of embryogenesis (approximately 66.7% of development at 28°C) (Figure 1D-J) were washed in distilled water, shock frozen in liquid nitrogen and stored at -80°C.

Embryos were collected from cages containing 25-50 females per cage. Ovaries from one adult female (Fig. 1B, C) were dissected from the body cavity, rinsed in 1X PBS, and homogenized in TRIzol (Invitrogen, NY, USA).

cDNA Synthesis

Total RNA was isolated separately from embryos at each day of embryonic development and from ovaries, using TRIzol (Invitrogen, NY, USA) and following manufacturer's instructions. RNA isolation was performed separately from embryonic and ovarian tissues, so that tissue lysis, which can affect the efficiency of subsequent RNA isolation, would be as homogeneous as possible within a sample. A pilot study was first conducted to determine library quality by sequencing ovarian and embryonic cDNA separately. For this pilot sequencing run, cDNA was synthesized using the SMART cDNA synthesis kit (Clontech, CA, USA) and normalized using the Evrogen Trimmer Direct kit (Evrogen, Moscow, Russia) following previously described methods [11]. Results from both libraries were comparable in read length and sequence quality, and all further experiments were carried out with pooled RNA libraries as described below. Raw reads from the pilot studies were incorporated into the final assembly as previously described [11].

To create a pooled cDNA library for large-scale sequencing, 1.5 µg of each of the mixed-stage embryonic RNA pool and ovarian RNA was used as a template for first strand cDNA synthesis. cDNA was synthesized as previously described [11]. Primary amplification proceeded with 10 PCR cycles monitored in real-time via qPCR [22], and secondary amplification began to plateau after 9 cycles. 16 parallel reactions of 0.73 µg each were co-purified into elution buffer using QIAquick PCR purification columns (Qiagen Inc., CA, USA). These 16 parallel reactions were identical, and were performed in individual tubes for the sole reason that a single PCR reaction sufficient to generate the 2 µg of cDNA required for sequencing would have had to be performed in a volume too large to undergo efficient cycling in our PCR machine (Bio-Rad Tetrad 2). We therefore calculated the predicted yield from the largest single PCR reaction that we could perform in our machine, and scaled up the number of reactions in parallel to achieve the required 2 µg total yield.

454 Titanium Pyrosequencing

The samples were nebulized, adaptor-ligated, and pyrosequenced using the 454 GS-FLX platform on pilot embryonic and ovarian cDNA separately, or the 454 GS-FLX Titanium platform for pooled ovarian/embryonic cDNA samples by the Institute for Genome Science and Policy DNA Sequencing Facility (Duke University). All of the raw reads generated in this study have been submitted to the NCBI Short Read Archive (Study Accession Numbers SRX023831, SRX023830, and SRX023832).

Sequence Assembly

Sequences were trimmed and assembled with Newbler v2.5, which was shown to outperform other assemblers for *de novo* assembly of 454 pyrosequencing reads [47]. Assembly parameters are described in [48], with the exception of the file used for the `-vt` flag (“Gb Adaptors”), which is available at <http://www.extavourlab.com/protocols/index.html>. Assembly results are available at <http://www.extavourlab.com/resources/index.html> and at <http://asgard.rc.fas.harvard.edu/download.html>.

Sequence Annotation

A nucleotide BLAST database was created using the isotigs and singletons produced by the Newbler assembly. To increase efficiency of BLAST comparison to this database, we first removed redundant isotigs and singletons created due to a combination of putative SNPs, sequencing errors, and low quality reads. Note that these data could in principle yield SNP data, but as we did not use an isogenic *G. bimaculatus* culture, nor do we have estimates of polymorphism for the culture, an accurate SNP analysis is not performed in the present study. Each assembly product was compared with the BLAST database using the BLASTN algorithm. Individual isotigs and singletons with BLAST hits (>95% identity based on bit score and sequence length) to longer sequences in the assembly, resulting in a high scoring segment pair (HSP) that spans the full length of the sequence, were removed. To identify the number of unique BLAST hits we followed the method described in [48].

To identify members of signaling pathways as described by the KEGG database [49], we manually annotated the *G. bimaculatus* transcriptome as described in [48].

Briefly, BLAST was used to compare the sequences of *D. melanogaster* pathway members with the *G. bimaculatus* transcriptome assembly and the top hit was selected as a putative ortholog with an E-value cutoff of e-10.

To determine whether the *de novo* assembly contained members of previously known *G. bimaculatus* GenBank accessions, we used tBLASTn (for 80 protein coding genes) or BLASTn (for 3 ribosomal RNA genes) to query the *G. bimaculatus* transcriptome assembly.

For automatic annotation of all transcriptome sequences, we designed a custom script called “Gene Predictor” (genePrediction.pl, available at <http://www.extavourlab.com/protocols/index.html>). This script assigns putative gene orthology based on comparisons with the *D. melanogaster* proteome, downloaded as described in Table S1. A protein BLAST database was created using the *D. melanogaster* proteome. A nucleotide BLAST database was created using the non-redundant assembly products (isotigs and singletons) of the *G. bimaculatus de novo* transcriptome assembly. The top 50 BLAST hits for each sequence of the *D. melanogaster* proteome compared with the *G. bimaculatus* transcriptome were obtained using the TBLASTN algorithm and stored in a MySQL database. Reciprocally, the top BLAST hit for each sequence of the *G. bimaculatus* transcriptome against the *D. melanogaster* proteome was obtained using the BLASTX algorithm and stored within a separate MySQL database. A custom script then iterates through each of the entries of the *D. melanogaster* proteome vs. the *G. bimaculatus* transcriptome MySQL database indices based on query identity and e-value. The same script also checks the *G. bimaculatus* transcriptome sequence identity against the *D. melanogaster* proteome

MySQL database to confirm if the reciprocal top BLAST hit is the same as the *D. melanogaster* query. After confirmation of the reciprocal BLAST identity, the script verifies whether any *G. bimaculatus* transcriptome sequences have already been assigned to the same *D. melanogaster* protein. If the existing sequence does not overlap with the confirmed sequence for more than 14 amino acids based on their HSP against the *D. melanogaster* protein, both sequences are recorded as orthologs. Otherwise, the confirmed sequence is further processed to determine whether it is a putative isoform or paralog of the existing sequence. If the confirmed sequence is a singleton or in the same isogroup as the existing sequence based on Newbler prediction, it is designated as an alternate isoform; otherwise, the sequence is annotated as a putative paralog.

A list of all curated *D. melanogaster* transcription factors was downloaded on March 26th 2011 from <http://flytf.org>. Each *D. melanogaster* transcription factor was examined to determine whether it was predicted to have an ortholog in the *G. bimaculatus* transcriptome using the Gene Predictor script described above. Custom scripts to generate tables based on the ASGARD schema (“ASGARD_NEW_DB.pl”) [50], upload assembled transcriptome sequences into ASGARD tables (“ASGARD_UPLOAD.pl”), upload BLAST results of the *D. melanogaster* proteome against the assembled transcriptome (“up_DMP.pl”), upload the BLAST results of the assembled transcriptome against the *D. melanogaster* proteome (“up_vDMP.pl”), and determine the best reciprocal BLAST result for each assembly products (“gene_prediction.pl”) are available at http://www.extavourlab.com/protocols/bio_tools/ASGARD_upload+Gene_Predictor.zip

.

Determination of sequencing depth and transcript completion

Ortholog hit ratio calculations and subassembly experiments were performed as described in [48]. Briefly, ortholog hit ratios were calculated using a custom script

("OrthologHitRatio.pl" available at

http://www.extavourlab.com/protocols/bio_tools/Perl_Transcriptome_Analysis_Scripts.zip

) that compares the length of each assembly product with the full length of its putative orthologous mRNA in *D. melanogaster*, based on the reciprocal best BLAST hit criteria described above. Subassemblies were performed by assembling progressively larger random subsets of all trimmed reads, using the same assembly parameters as those used for the complete assembly.

Protein Domain Analysis

23 proteomes based on completely sequenced genomes and two EST libraries were downloaded as described in Table S1. A protein BLAST database was created from each proteome. All *G. bimaculatus* assembly products were compared with each database using the BLASTX algorithm with an E-value cutoff of 1e-5. The resulting reports were parsed using the Uniqueblast.pl script as previously described [48] (available at

<http://www.extavourlab.com/protocols/index.html>).

A local installation of EST Scan [51] (ESTSCAN 3.03) was downloaded on April 11th 2011 as a Linux rpm package from <http://estscan.sourceforge.net/>. All assembly products were screened using ESTSCAN with default parameters, except for the "-l" flag that was used with a value of 20 to restrict the minimum result size to 20 amino acids.

The “-t” flag was also used to allow ESTSCAN to produce the predicted protein sequence of each assembly product.

A local installation of InterPro Scan [52,53] (IPRSCAN 4.7) was downloaded on April 15th 2011 from <ftp://ftp.hgc.jp/pub/mirror/ebi/software/iprscan/index.html>. The “-cli” flag was used to turn on pipeline mode and suppress html outputs. All assembly products were screened using IPRSCAN against existing protein feature databases [54,55,56,57,58,59,60,61,62,63,64,65], and the results were stored in xml format for further analysis.

Welch’s t-test (appropriate in this case for use with samples with unequal variance [66]) was used for statistical comparisons of lengths of sequences and predicted protein coding regions in various annotation categories.

Results and Discussion

Collection and preparation of material

We aimed to create a transcriptome containing genes deployed during oogenesis, when maternally deposited factors required for embryogenesis may be synthesized, and during all stages of embryogenesis. We therefore collected ovaries (Figure 1B, C) and embryos from early to late stages of embryogenesis (Figure 1D-J) for mRNA extraction. We pooled these mRNA samples and prepared non-normalized cDNA libraries for 454 Titanium pyrosequencing. We chose to omit normalization in preparing these libraries as our previous studies [11] suggest that at this scale of sequencing, normalization does not significantly aid in gene discovery.

Sequencing and basic transcriptome assembly

We used Newbler v2.5 (Roche) for the *de novo* assembly of 4,248,348 raw reads (1,483,726,666 bp) obtained by 454 pyrosequencing (Table 1). Using default Newbler assembly parameters, raw reads were screened and trimmed of both 5' and 3' adaptors (see Methods), and low quality reads were removed. (Newbler's quality scores are defined as "Phred-like" or "Phred equivalent" [67]. The Phred quality score is a widely used base quality parameter defined by determining qualities of the data used to generate each base call [68,69]. We used a Newbler quality score cutoff of >20; a Phred score of 20 would indicate a base call accuracy of $\geq 99\%$). 99.26% of all reads passed this quality control process (4,216,721 reads = 1,449,059,795 bp) (Figure S1A, Table 1), and were subsequently used in the sequence alignment process. 88.78% of these reads (3,743,561) were fully assembled, meaning that the entire read sequence was used in a contig. 6.69% (282,259) were partially assembled, meaning that the entire read was not used in a contig (Figure S1B, C). Of the 190,901 good quality reads (4.53%) that were not aligned, 13,416 (0.32%) were too short (<40 bp) to be included in the assembly, 1,989 (0.05%) were predicted to be from a repeat region (meaning that >70% of the read's seeds match at least 70 other reads, or determined to partially overlap a contig; note that portions of reads in this category that overlap unique contigs are still included in the assembly results), 54,691 (1.30%) were considered outliers (e.g. chimeric reads or results of sequencing errors), and 120,805 (2.86%) were preserved as singletons.

Newbler assembly products fall into one of four categories: (1) *contigs* are groups of assembled reads with significant overlapping regions (we used the Newbler default

minimum overlap of 40 bp), which may represent exons; (2) *isotigs* are continuous paths through a given set of contigs, and represent putative transcripts, including possible splice variants of a given transcription unit; (3) *isogroups* are groups of isotigs that were assembled from the same contig set, and are the closest to gene predictions as it is possible for a *de novo* assembly to achieve; and (4) *singletons*, which are single good quality reads that lack significant overlap with any other read, and therefore are not incorporated into any contig. We use these terms henceforth to refer to the *G.*

bimaculatus assembly products. It is important to note that determination of whether contigs represent true exons, or isotigs true transcripts, would require further validation by sequencing full-length cDNAs and comparison with a fully sequenced genome. For this reason we refer to the *G. bimaculatus* transcriptome *de novo* assembly products as “contigs” and “isotigs” or “predicted transcripts” or “putative transcripts” throughout, rather than as “exons” or “transcripts” respectively.

Upon assembly we obtained 43,321 unique contigs using the aligned reads (Table 1). Newbler then further assembled these contigs into 21,512 isotigs that belonged to 16,456 isogroups (Table 2). 13,157 (79.95%) of the isogroups (putative genes) consist of only a single isotig, and on average there are 1.2 isotigs per isogroup (Table 2). 12,701 (62.78%) isotigs consist of a single contig, and on average there are 1.7 contigs per isotig. The isotig N50 is 2,133 bp (Table 1), meaning that the majority of predicted transcripts are over 2kb in length. FASTA files of all assembly products are available for download from our interactive database (described below).

Assessment of transcript coverage and depth

The average coverage across the assembly is 51.3 reads per base pair; in other words, each base pair of the assembly was sequenced on average over 50 times. This coverage is high compared to other *de novo* transcriptome assemblies [11,48,70], which we attribute largely to the high number of reads used to create the *G. bimaculatus* transcriptome. We note, however, that the *G. bimaculatus* transcriptome coverage we obtained is more than twice as high as that of the recently *de novo* assembled transcriptome for the crustacean *Parhyale hawaiiensis* (25.4 reads/bp), even though the *G. bimaculatus* transcriptome contained only 1.3 fold more base pairs in raw reads than that of *P. hawaiiensis*, which was also generated from embryonic and ovarian cDNA, and was assembled and annotated identically to the *G. bimaculatus* transcriptome described in this report [48].

An additional measure of coverage is the average contig read depth (total number of base pairs from all reads aligned to generate a given contig, divided by contig length). This value is 391 bp/contig, with a median value of 16.7 bp/contig. We note that the predicted transcript coverage (number of base pairs of raw reads comprising each contig) is highly variable, suggesting that some genes are represented by many more raw reads than others (Figure 2). 19,093 (43.97%) contigs had a coverage ≤ 10 bp/contig, and 538 contigs (1.24%) had a coverage $\geq 10,000$ bp/contig.

We wished to determine whether similar coverage levels and predicted transcript lengths could have been obtained with fewer reads, and how well our transcriptome had identified all putative transcripts present in our samples. To do this, we created subassemblies using randomly chosen subsets of reads, starting with 10% of reads and adding increments of 10% up to the full complement of trimmed reads. For each subset of reads, we performed an independent assembly with Newbler v2.5. For each of these nine

subassemblies, we then assessed both read length distribution and the number of unique BLAST hits against the NCBI non-redundant protein database (**nr**) with an E-value cutoff of 1e-10. The mean coverage per bp was strongly positively correlated ($R^2 = 0.96$, linear regression) with the number of reads used for the assembly (Figure 3A, blue line). We also found that as the number of reads used in the subassembly increased, the proportion of reads left as singletons decreased from 11.25% for the 10% subassembly, to 2.86% in the full assembly. This is likely because contigs and isotigs increased in length as reads were added (Figure 3B), as we observed an increase in isotig N50 from 1,290 bp with 10% of reads to 2,133 bp with all reads. The distribution of isotig lengths in each subassembly (Figure 3B) indicates the maximum length of assembled isotigs given a certain number of reads. A small proportion of isotigs exceeding 4 kb can be obtained with only 10% of all reads, but by assembling all reads it was possible to obtain predicted transcripts exceeding 10 kb (Figure 3C).

The number of unique BLAST hits against **nr** obtained from all isotigs also increased with the number of reads (Figure 3A, red line), but at a slower rate than that of mean coverage per bp (Figure 3A, blue line). Slightly fewer unique BLAST hits were obtained from isotigs generated with 100% of reads compared to 90%, which may mean that previously unconnected contigs were increasingly incorporated into isotigs as they increased in length and acquired overlapping regions.

To estimate the degree to which full-length transcripts might be predicted by the transcriptome, we determined the ortholog hit ratio [70] of all assembly products by comparing the BLAST results of the full assembly against the *Drosophila melanogaster* proteome. The ortholog-hit ratio is calculated as the ratio of the length of a transcriptome

assembly product (isotig or singleton) and the full length of the corresponding transcript. Thus, a transcriptome sequence with an ortholog hit ratio of 1 would represent a full-length transcript. In the absence of a sequenced *G. bimaculatus* genome, for the purposes of this analysis we use the length of the cDNA of the best reciprocal BLAST hit against the *D. melanogaster* proteome as a proxy for the length of the corresponding transcript. For this reason, we do not claim that an ortholog hit ratio value indicates the true proportion of a full-length transcript, but rather that it is likely to do so. The full range of ortholog hit ratio values for isotigs and singletons is shown in Figure 4. Here we summarize two ortholog hit ratio parameters for both isotigs and singletons: the proportion of sequences with an ortholog hit ratio ≥ 0.5 , and the proportion of sequences with an ortholog hit ratio ≥ 0.8 . We found that 63.8% of *G. bimaculatus* isotigs likely represented at least 50% of putative full-length transcripts, and 40.0% of isotigs were likely at least 80% full length (Figure 4B). For singletons, 6.3% appeared to represent at least 50% of the predicted full-length transcript, and 0.9% were likely at least 80% full length (Figure 4B). Most ortholog hit ratio values were higher than those obtained for the *de novo* transcriptome assembly of another hemimetabolous insect, the milkweed bug *Oncopeltus fasciatus* [11] (Figure 4A, B). We suggest that this may be explained by the fact that the *G. bimaculatus de novo* transcriptome assembly contains transcript predictions of higher coverage and longer isotigs (N50 = 2,133 compared to 1,735 for *O. fasciatus* [11]) that are likely closer to predicted full-length transcript sequences, relative to the *O. fasciatus de novo* transcriptome assembly [11]. However, we cannot exclude the possibility that the higher ortholog hit ratios obtained with the *G. bimaculatus* transcriptome may be due to its greater sequence similarity with *D. melanogaster* relative

to *O. fasciatus*. Genome sequences for the two hemimetabolous insects, and rigorous phylogenetic analysis for each predicted gene in both transcriptomes, would be necessary to resolve the origin of the ortholog hit ratio differences that we report here.

Annotation using BLAST against the NCBI non-redundant protein database

All assembly products were compared with the NCBI non-redundant protein database (**nr**) using BLASTX. We found that 11,943 isotigs (55.52%) and 10,815 singletons (8.95%) were similar to at least one **nr** sequence with an E-value cutoff of 1e-5 (henceforth called “significant similarity”). The total number of unique BLAST hits against **nr** for all non-redundant assembly products (isotigs + singletons) was 19,874, which could correspond to the number of unique *G. bimaculatus* transcripts contained in our sample. The *G. bimaculatus* transcriptome contains more predicted transcripts than other orthopteran transcriptome projects to date (Table 1). This may be due to the high number of bp incorporated into our *de novo* assembly, which was generated from approximately two orders of magnitude more reads than previous Sanger-based orthopteran EST projects [71,72,73,74,75]. However, we note that even a recent Illumina-based locust transcriptome project that assembled over ten times as many base pairs as the *G. bimaculatus* transcriptome, predicted only 11,490 unique BLAST hits against **nr** [71]. This may be because the tissues we samples possessed a greater diversity of gene expression than those for the locust project, in which over 75% of the cDNA sequenced was obtained from a single nymphal stage [71]. Although we have used the *de novo* assembly method that was recommended as outperforming other assemblers in

analysis of 454 pyrosequencing data [47], we cannot exclude the possibility that under-assembly of our transcriptome contributes to the high number of predicted transcripts

Since isogroups are groups of isotigs that are assembled from the same group of contigs, the isogroup number of 16,456 may represent the number of *G. bimaculatus* unique genes represented in the transcriptome. However, because by definition *de novo* assemblies cannot be compared with a sequenced genome, several issues limit our ability to estimate an accurate transcript or gene number for *G. bimaculatus* from these ovary and embryo transcriptome data alone.

The number of unique BLAST hits against **nr** (19,874) or isogroups (16,456) may overestimate the number of unique genes in our samples, because the assembly is likely to contain sequences derived from the same transcript but too far apart to share overlapping sequence; such sequences could not be assembled together into a single isotig and would therefore have been considered “different genes.” If such assembly products were derived from different regions of the same transcript and obtained distinct BLAST hits against **nr**, then these would be counted as two unique BLAST hits against **nr**. This limitation is an inevitable result of performing *de novo* assembly in the absence of a reference genome, and is unavoidable in the case of *G. bimaculatus* as no orthopteran genomes have yet been sequenced. Conversely, the number of unique BLAST hits against **nr** could underestimate the number of unique genes, because they cannot include those isotigs (9,569 = 44.5% of all isotigs) and singletons (109,990 = 91.0% of all singletons) that lacked significant BLAST hits against **nr**. Such sequences could represent non-coding sequences with no matches to the coding-region data contained in **nr**, or could lack sufficient similarity to known sequences. Finally, because

our transcriptome libraries were prepared only from ovarian and embryonic tissue, it is unlikely to contain transcripts of all *G. bimaculatus* genes, many of which could be expressed exclusively postembryonically and/or in specific nymphal or adult tissue types. Determination of the total gene number for *G. bimaculatus* must therefore await complete genome sequencing.

We wished to understand the relative similarities of the *G. bimaculatus* transcriptome sequences to those from other organisms. Specifically, we asked what proportion of genes found in sequenced animal genomes had putative orthologs in the *G. bimaculatus* transcriptome. To this end, we used BLAST to compare each non-redundant assembly product (E-value cutoff 1e-5) to the proteomes of several organisms with completely sequenced genomes (Table S1). We found that overall, 33.49% of the sequences contained in insect proteomes had matches in the *G. bimaculatus de novo* transcriptome assembly, compared to 22.28% of sequences from deuterostome proteomes (Figure 5). Within the insects, the proportion of hits to the *D. melanogaster* proteome was lower than the proportion of hits to most other insects. This may reflect the relatively greater divergence from a last common insect ancestor, as *D. melanogaster* belongs to the most derived insect order, the Diptera. However, we noted that the proportion of matches to some insect proteomes appeared unusually low given their phylogenetic relationship to Orthoptera. Specifically, only 18.1% of proteome sequences from the aphid *Acyrtosiphon pisum*, a hemimetabolous insect, had hits in the *G. bimaculatus* transcriptome, compared with an average of 36.1% across all holometabolous proteomes surveyed (Figure 5). This is consistent with the description of the *A. pisum* genome containing many unusual features relative to other insect genomes, including extensive

gene family duplications and gene loss [6,76,77,78]. The relatively high proportion of holometabolous proteome sequences with matches in the *G. bimaculatus* transcriptome suggests that these organisms may share more features derived from a last common insect ancestor than does *A. pisum*, and highlights the need for further genomic resources in the Hemimetabola. We caution that there are limitations to the biological information that can be derived from these comparisons, as not all animal genomes used for this analysis have comparable levels of coverage or annotation.

Manual annotation of conserved developmental genes and members of signaling pathways

G. bimaculatus has been the subject of molecular embryology for over a decade, and as a result over 80 GenBank accessions are available (NCBI accessed 12 August 2012). We asked whether these genes were represented in our transcriptome, and found that 72.3% of them were present (60/83). Moreover, the transcriptome contributed to these accessions by extending their sequences by an average of 737 nucleotides per accession (205.0% on average across all 83 *G. bimaculatus* GenBank accessions) and in some cases by over 1,700% (Table S2). This shows that the *G. bimaculatus* transcriptome will be an extremely useful resource for continued research into the function and evolution of most previously cloned genes.

To determine the transcriptome's utility as a source of new gene discovery, we searched for putative orthologs of the 1,168 *D. melanogaster* transcription factors catalogued in the FlyTF transcription factor database [79]. We found that 542 (46.4%) of them were present, based on the criterion of being the best reciprocal BLAST hit with a

D. melanogaster sequence using an E-value cutoff of 1e-5 (Table S3). We also undertook manual annotation of 122 genes from seven conserved metazoan signaling pathways (Table S4), 261 genes involved in male and female gametogenesis in *D. melanogaster* (Table S5), and 24 additional genes with roles in maternal or zygotic embryonic patterning (Table S6). For the Notch [80], TGF-beta [81], Wnt [82], JAK/STAT [83], MAPK [84] and *hedgehog* [85] signaling pathways, most *G. bimaculatus* orthologs of these genes were previously unknown. Our transcriptome newly identified 66 genes participating in these signaling pathways (Table S4, Figure S2), including nearly all members besides the ligand of the *hedgehog* pathway (Figure 6A). In the case of the Hippo signaling pathway [86], for which most *G. bimaculatus* core kinase orthologs were already present in GenBank, the *G. bimaculatus de novo* transcriptome assembly increased the length of known sequences by an average of 323%, and by as much as 1,119% in the case of the *discs overgrown (dco)* gene (Figure 6B, Table S2).

Automated annotation using the custom script “Gene Predictor” identifies 14,130 transcriptome sequences as putatively orthologous to D. melanogaster genes

Although manual annotation proved a highly effective way to identify developmental genes of interest in the *G. bimaculatus* transcriptome, it is not efficient at large scales. We therefore developed an automated annotation tool that uses the criterion of best reciprocal BLAST hit against the *D. melanogaster* proteome (E-value cutoff 1e-5) to propose putative orthologs for all assembly products of the transcriptome. This method is not qualitatively different from manual annotation using BLAST with a specific known

sequence as a query, but rather simply automates the process of detecting a best reciprocal BLAST hit, which is a method of orthology assignment routinely employed as an annotation method in genomics studies using insect genomes [87,88,89]. Using this tool, called Gene Predictor (see Methods), we were able to assign putative orthologs to 43.7% of isotigs, very close to the proportion of isotigs (55.5%) with significant BLAST hits against **nr** (Figure 7A). Of the 60 known *G. bimaculatus* GenBank accessions that were identified in the transcriptome by manual annotation (Table S2), 52 have significant BLAST hits to a *D. melanogaster* gene (the remaining 8 genes have significant similarity only to non *D. melanogaster* sequences, as determined by BLAST against **nr**). Gene Predictor correctly identified 36 of these 52 genes (69.2%). Gene Predictor's failure to identify the remaining 16 genes (30.8%) means that while these genes do have significant BLAST hits in the *D. melanogaster* genome, they are more similar to a non-*D. melanogaster* gene, and are thus not the reciprocal best BLAST hit of any *D. melanogaster* gene.

These results suggest that for *de novo* insect transcriptome assemblies, Gene Predictor could be an efficient annotation tool, as it is nearly as effective as BLAST mapping against the large **nr** database, but is computationally much less intensive as it relies only on the *D. melanogaster* proteome of 23,361 predicted proteins. Relative to BLAST mapping against **nr**, Gene Predictor was more effective at suggesting orthologs for isotigs than for singletons (Figure 7A), likely due to the fact that isotigs are easier to map by any method as they contain more sequence data. Gene Predictor did not, however, assign orthologs to any assembly products that did not already have a significant BLAST hit in **nr** (Figure 7B), as expected since the *D. melanogaster*

proteome is contained within **nr**. Conversely, not all assembly sequences with BLAST hits in **nr** obtained a significant hit with Gene Predictor (Figure 7B), indicating that some of the *G. bimaculatus* predicted transcripts share greater similarity to sequences other than those in the *D. melanogaster* proteome, or may represent genes that have been lost in *D. melanogaster*. The Gene Predictor scripts are freely available at <http://www.extavourlab.com/protocols/index.html>.

Transcripts lacking significant BLAST hits against nr may encode functional protein domains

The majority (55.5%) of predicted transcripts retrieved a significant BLAST hit against the **nr** database (Figure 7A). This exceeds the proportion of *de novo* assembly products typically identifiable by BLAST mapping against **nr** [70], including the 43.4% and 29.5% of predicted transcripts mapped in this way from two *de novo* arthropod transcriptome assemblies that we previously constructed using similar methods to those described here [11,48]. This may be due to the much higher read depth and coverage of the *G. bimaculatus* transcriptome, which to our knowledge is the largest *de novo* assembled transcriptome available for the Hemimetabola, and the largest 454-based transcriptome for any organism to date. Even this assembly, however, contains a large proportion (44.5%) of sequences of unknown identity. These sequences could represent contaminants of unknown origin, sequences that are too short to obtain significant hits to **nr** sequences, non-coding transcripts, non-coding portions of protein-coding transcripts, or clade- or species-specific transcripts that may be unidentifiable due to the paucity of orthopteran genomic data in GenBank. We believe that significant contaminants are

unlikely, as less than one percent of all assembly products retrieved BLAST hits to prokaryote, fungal or plant sequences with an E-value cutoff of $1e-10$.

We also compared the length (in nucleotides) of sequences with and without significant BLAST hits (Tables 3, 4), and found that unidentified isotigs were significantly shorter than isotigs with BLAST hits (Table 5). The difference was also significant for singletons (Tables 4, 5). This is consistent with the possibility that contig length may play a role in sequence recognizability, also observed with the low proportion of singletons with significant BLAST hits compared to isotigs (9.0% vs 55.5%; Figure 8A, B).

To obtain additional biological information about sequences that failed to obtain significant BLAST hits against **nr**, we therefore applied EST Scan analysis to determine whether these sequences potentially encoded unknown proteins. EST Scan uses known differences in hexanucleotide usage between coding and non-coding regions to detect potential coding regions in DNA sequences, without requiring open reading frames [51]. We found that 2,468 (25.8%) unidentified isotigs and 16,409 (14.9%) unidentified singletons were predicted to contain protein-coding regions (Figure 8). Isotigs without predicted coding regions were significantly shorter than sequences with predicted coding regions (Tables 3, 5); the difference was also significant for singletons (Tables 4, 5). Previously unidentified isotigs that were protein-coding were significantly shorter than isotigs with significant BLAST hits, and encoded significantly fewer amino acids (Tables 3, 5, 6). This may mean that significant BLAST hits were not obtained for some of these sequences either because of insufficient contig lengths, or because they contained relatively less protein-coding content, or both. These observations demonstrate that

although these 18,877 sequences are not significantly similar to known proteins in the NCBI nr database, they may nevertheless represent portions of coding rather than non-coding transcripts.

We then used InterPro Scan [52,53] to query predicted coding regions for predicted functional protein domains. InterPro Scan queries the InterPro consortium databases (ProDom [54], PRINTS [90], SMART [56], TIGRFAMs [57], Pfam [58], PROSITE [59], PIRSF [60], SUPERFAMILY [61], CATH [62], PANTHER [63], SignalPHMM [64], and Transmembrane [65]) for signatures of protein domains of known function. It also identifies evolutionarily conserved protein domains that are predicted to be functional based on their conservation but have no described molecular function to date, called Domains of Unknown Function (DUFs) [91]. This analysis revealed that of those protein-coding sequences of unknown identity, 495 (20.0%) isotigs and 1,447 (6.7%) singletons were predicted to contain functional protein domains. These results show that 1,942 sequences from the *de novo* transcriptome assembly that could not be identified based on BLAST against **nr** alone may nonetheless encode functional proteins present during *G. bimaculatus* oogenesis and embryogenesis.

*Taxonomic bias of the **nr** database can limit gene identification in de novo assembled transcriptomes*

Because orthopteran sequence data are poorly represented in **nr**, we asked whether at least some of the *G. bimaculatus* transcriptome sequences that appeared to lack significant similarity to known genes might show similarity to sequences from other orthopterans available in the form of EST collections. To determine this, we compared

the 9,569 isotigs (44.5% of all isotigs) and 109,990 singletons (91.0% of all singletons) from the *G. bimaculatus* transcriptome that lacked significant **nr** hits, with the EST collections for the orthopterans *L. migratoria* and *L. kohalensis*. *L. migratoria* of the suborder Caelifera (grasshoppers and locusts) is a migratory locust that is widespread throughout Asia, Africa, and Australasia [92], and is heavily studied due to its impact as an agricultural pest [e.g. 93,94]. The available sequence collections for this locust sampled transcripts from larval stages L4 and L5 [71,72,73], which is when transition between the solitary and gregarious (swarming) behavior of these locusts becomes irreversible [73,95]. *L. kohalensis* belongs to the suborder Ensifera (katydids and crickets), and is a Hawaiian species that has been used extensively for studies of the physiology and evolution of speciation and acoustic preference [e.g. 23,96,97]. The EST library available for this cricket contains sequences derived from transcripts of the larval central nervous system [74]. Because these data are derived from EST collections, they are available through GenBank but are not included in **nr**.

Using BLAST with an E-value cutoff of e^{-5} , we found that the majority of previously unidentified *G. bimaculatus* transcriptome sequences also lacked significant similarity to *L. migratoria* or *L. kohalensis* sequences. This may be due to the difference in starting material for the libraries compared, as the *G. bimaculatus* transcriptome contains transcripts from ovaries and embryos, while the other two libraries represent exclusively post-embryonic transcripts, and the *L. kohalensis* library is further restricted only to transcripts from the nervous system. However, 406 isotigs (4.24%) and 1,058 singletons (0.96%) did display significant similarity (Figure 9A, B), suggesting that these transcripts could represent “orthopteroid” genes. However, we noted that sequences of

both isotigs and singletons in this category contained significantly fewer nucleotides (Tables 3-5) and encoded significantly fewer amino acids on average (Tables 3, 4, 6) than transcriptome sequences with BLAST hits to **nr** (Tables 3-6). An alternative explanation for these apparent “orthopteroid” sequences is thus that these sequences, as well as their matches from *L. kohalensis* and *L. migratoria*, might prove significantly similar to other sequences from **nr**, if their transcript sequences were longer.

Because Ensifera and Caelifera are believed to have diverged 300 Mya [5], we predicted that we would find greater similarity between sequences from the two crickets, than between *G. bimaculatus* and the locust. Accordingly, of the putative “orthopteroid sequences,” 746 (51.0%) *G. bimaculatus* sequences yielded hits exclusively to *L. kohalensis* sequences, compared to 156 (10.7%) sequences with exclusive hits among *L. migratoria* sequences (Figure 9C'). This likely reflects the closer phylogenetic relationship between the two crickets, which are both within the same family of Gryllidae.

Putative orthopteroid-specific sequences contain a high proportion of predicted protein coding domains of unknown function (DUFs)

Finally, we asked whether these “orthopteroid sequences” shared any characteristics that might aid in understanding their putative clade-specific functions. We used InterPro Scan [52] to determine the distribution of recognizable protein domains among transcriptome sequences with significant *L. kohalensis* or *L. migratoria* hits, and compared them with those of all transcriptome sequences with significant BLAST hits to **nr**. We found that the number of distinct domains was similar for *L. kohalensis*-like

sequences (77 different protein domains) and all other transcriptome sequences with significant BLAST hits (83 different protein domains), but considerably lower for *L. migratoria*-like sequences (55 different protein domains). Given the small number of sequences examined here (Figure 9C), this is unlikely to represent true differences in protein type between the three datasets.

However, the datasets differed strikingly in the relative proportions of different protein domains encoded. Considering the top 25 most frequently represented protein domains within each dataset, the most abundant domains in both orthopteran-like groups were domains of unknown function (DUFs, 18.8% of both orthopteran matches combined), followed by ubiquitin family domains (Pfam PF00240, 10.9%), zinc finger domains (multiple Pfam categories combined, 10.2%), and RNA recognition motifs (Pfam PF00076, 5.5%) (Figure 10A, B). In contrast, transcriptome sequences with significant BLAST hits to **nr** encoded proteins principally containing zinc finger domains (multiple Pfam categories combined, 22.7%), protein kinase domains (Pfam 00069, 16.2%), and ankyrin repeat domains (Pfam PF00023, 12.0%), followed by RNA recognition motifs (Pfam PF00076, 9.6%) and BTB/POZ domains (Pfam PF00651, 9.0%) (Figure 10C). These differing proportions of predicted protein domains between orthopteran-matched and **nr**-matched *G. bimaculatus* sequences were observed even when all predicted protein domains were considered (Figure S3). We speculate that the “orthopteroid-like” proteins predicted to be present in the *G. bimaculatus* transcriptome might share greater functional similarity with orthopteran proteins than with proteins from other organisms represented in **nr**. Moreover, the high proportion of DUFs predicted in these “orthopteroid-like” proteins may mean that some of these DUFs serve

clade-specific functions. The specific roles of these genes in *G. bimaculatus* and other orthopterans are currently unknown, and will require functional genetic testing to be elucidated. However, the present analysis demonstrates that even for *de novo* assembled transcriptome sequences that are not easily identifiable based on GenBank comparisons, it may be possible to extract potentially meaningful biological and evolutionary information, and with further refinement, perhaps even to define new or clade-specific DUFs as candidates for future functional testing.

Creation of a searchable database to house arthropod de novo assembled transcriptomes

The volume of high-throughput transcriptome data available for all organisms is rapidly increasing, but many of these datasets are not publicly available in an easily searchable format. The NCBI Short Read Archive [98] provides a repository for raw read data from transcriptome projects, but a searchable interface for *de novo* assembled transcriptomes that do not have an associated genome sequence or previously developed community web interface is lacking. Like EST collections, transcriptome assemblies can be made public through the NCBI Transcriptome Shotgun Assembly Sequence Database (TSA: <http://www.ncbi.nlm.nih.gov/genbank/tsa>), but annotation of these data is not required, and they are not included in **nr**. To maximize the public utility of our data, we therefore created a searchable database that facilitates access to the annotated *G. bimaculatus de novo* assembled transcriptome reported here. The Assembled Searchable Giant Athropod Read Database (ASGARD) includes all **nr** BLAST, manual annotation, and Gene Predictor annotation results for the *G. bimaculatus* transcriptome. Details of the design and database schema of ASGARD have been previously described [50]. This

database also contains two additional *de novo* assembled transcriptomes that we constructed previously, for the milkweed bug *Oncopeltus fasciatus* [11] and the amphipod crustacean *Parhyale hawaiiensis* [48]. The *O. fasciatus* transcriptome, which was originally assembled with Newbler v2.3 [11], was re-assembled with Newbler 2.5, which was used to assemble the *P. hawaiiensis* and *G. bimaculatus* transcriptomes. Complete updated assembly files in FASTA format for all three transcriptomes can be downloaded via ASGARD. We also processed the *O. fasciatus* and *P. hawaiiensis* transcriptomes with the EST Scan, InterPro Scan, and the Gene Predictor script, so that they could be searched in the same way as the *G. bimaculatus* transcriptome. ASGARD allows users to search these *de novo* assembled transcriptomes in four ways: (1) for putative orthologs to known *D. melanogaster* genes (based on Gene Predictor results); (2) by searching the text of the top 50 significant BLAST hits for the name of any gene of interest (based on **nr** BLAST mapping results); (3) by searching for transcripts with a given GO term assignment; and (4) by read name if the unique identifier of a given assembly product is known (this information is provided in the results of the previous three searches). All search result output pages allow users to view and download the nucleotide sequences of matching assembly products, the pre-computed results of a BLAST search of that sequence against **nr** (E-value cutoff 1e-5), their predicted translation products if applicable (determined using EST Scan), and any predicted functional protein domains (determined using InterPro Scan). Finally, ASGARD also contains a BLAST interface that allows users to search any or all transcriptomes using the BLASTN, TBLASTN or TBLASTX algorithms. ASGARD is available at <http://asgard.rc.fas.harvard.edu>.

Acknowledgments

Thanks to Sumihare Noji for encouragement and funding contributions (JSPS KAKENHI 22124003/22370080), Evelyn E. Schwager and Franz Kainz for assistance with embryo collection, Lisa Bukovnik for administration of the sequencing, the Harvard Faculty of Arts and Sciences Research Computing group for database discussions, and the Extavour lab for discussions of the results. This work was partially supported by Harvard Stem Cell Institute Seed Grant SG-0057-10-00, Ellison Medical Foundation New Scholar Award AG-NS-07010-10, and NSF Grant IOS-0817678 to CE, an NSF Predoctoral Fellowship to BEC, a Fletcher Family award from Bowdoin College to Hadley Horsch, DFG Collaborative Research Centre 680 funds to Siegfried Roth, and JSPS KAKENHI 22124003/22370080/23687033 to Taro Mito.

Author Contributions

VZ performed experiments, helped design data analysis and analyzed data. BEC helped design research, performed experiments, collected and analyzed data. CGE, HH, SR, and TM obtained funding for the research. CGE proposed the idea for the research, helped design the research and analyze the data, wrote the manuscript with input from VZ and obtained funding for the research. All authors read and approved the final manuscript.

Competing Interests

The authors declare that they have no competing interests. Victor Zeng is currently employed by Stylux Incorporated. This does not alter our adherence to all the PLoS ONE policies on sharing data and materials.

References

1. Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, et al. (2002) The genome sequence of the malaria mosquito *Anopheles gambiae*. *Science* 298: 129-149.
2. Brown SJ, Denell R, Gibbs R, Klingler M, Lorenzen M, et al. (2008) The genome of the model beetle and pest *Tribolium castaneum*. *Nature* 452: 949-955.
3. Adams MD, Celniker SE, Holt RA, Evans CA, Gocayne JD, et al. (2000) The genome sequence of *Drosophila melanogaster*. *Science (New York, NY)* 287: 2185-2195.
4. Myers EW, Sutton GG, Delcher AL, Dew IM, Fasulo DP, et al. (2000) A whole-genome assembly of *Drosophila*. *Science (New York, NY)* 287: 2196-2204.
5. Grimaldi D, Engel MS (2005) *Evolution of the Insects*. Cambridge: Cambridge University Press. 772 p.
6. Consortium IAG (2010) Genome sequence of the pea aphid *Acyrtosiphon pisum*. *PLoS Biology* 8: e1000313.
7. Huebner E (2007) The *Rhodnius* Genome Project: The promises and challenges it affords in our understanding of reduviid biology and their role in Chagas' transmission. *Comparative Biochemistry and Physiology, Part A* 148: S130.
8. Megy K, Emrich SJ, Lawson D, Campbell D, Dialynas E, et al. (2012) VectorBase: improvements to a bioinformatics resource for invertebrate vector genomics. *Nucleic Acids Res* 40: D729-734.
9. Kirkness EF, Haas BJ, Sun W, Braig HR, Perotti MA, et al. (2010) Genome sequences of the human body louse and its primary endosymbiont provide insights into the permanent parasitic lifestyle. *Proceedings of the National Academy of Sciences* 107: 12168-12173.
10. Zdobnov EM, Bork P (2007) Quantification of insect genome divergence. *Trends in Genetics* 23: 16-20.
11. Ewen-Campen B, Shaner N, Panfilio KA, Suzuki Y, Roth S, et al. (2011) The maternal and embryonic transcriptome of the milkweed bug *Oncopeltus fasciatus*. *BMC Genomics* 12: 61.
12. Huber F, Moore TE, Loher W, editors (1989) *Cricket Behavior and Neurobiology*. Ithaca, NY: Cornell University Press. 571 p.
13. Hartfelder K, Emlen DJ (2011) Endocrine control of insect polyphenism. In: Gilbert LI, editor. *Insect Endocrinology*: Elsevier. pp. 464-522.
14. Zera AJ (2006) Evolutionary genetics of juvenile hormone and ecdysteroid regulation in *Gryllus*: a case study in the microevolution of endocrine regulation. *Comp Biochem Physiol A Mol Integr Physiol* 144: 365-379.
15. Tomioka K, Matsumoto A (2010) A comparative view of insect circadian clock systems. *Cellular and Molecular Life Sciences* 67: 1397-1406.
16. Danbara Y, Sakamoto T, Uryu O, Tomioka K (2010) RNA interference of *timeless* gene does not disrupt circadian locomotor rhythms in the cricket *Gryllus bimaculatus*. *Journal of Insect Physiology* 56: 1738-1745.
17. Tomioka K, Uryu O, Kamae Y, Umezaki Y, Yoshii T (2012) Peripheral circadian rhythms and their regulatory mechanism in insects and some other arthropods: a

- review. *Journal of Comparative Physiology B, Biochemical, Systemic, and Environmental Physiology* 182: 729-740.
18. Shaw KL, Parsons YM, Lesnick SC (2007) QTL analysis of a rapidly evolving speciation phenotype in the Hawaiian cricket *Laupala*. *Mol Ecol* 16: 2879-2892.
 19. Howard DJ, Marshall JL, Hampton DD, Britch SC, Draney ML, et al. (2002) The genetics of reproductive isolation: a retrospective and prospective look with comments on ground crickets. *American Naturalist* 159 Suppl 3: S8-S21.
 20. Ellison CK, Wiley C, Shaw KL (2011) The genetics of speciation: genes of small effect underlie sexual isolation in the Hawaiian cricket *Laupala*. *J Evol Biol* 24: 1110-1119.
 21. Maroja LS, Clark ME, Harrison RG (2008) Wolbachia plays no role in the one-way reproductive incompatibility between the hybridizing field crickets *Gryllus firmus* and *G. pennsylvanicus*. *Heredity* 101: 435-444.
 22. Andres JA, Maroja LS, Harrison RG (2008) Searching for candidate speciation genes using a proteomic approach: seminal proteins in field crickets. *Proceedings of the Royal Society of London Series B: Biological Sciences* 275: 1975-1983.
 23. Shaw KL, Lesnick SC (2009) Genomic linkage of male song and female acoustic preference QTL underlying a rapid species radiation. *Proceedings of the National Academy of Sciences of the United States of America* 106: 9737-9742.
 24. Fedorka KM, Mousseau TA (2004) Female mating bias results in conflicting sex-specific offspring fitness. *Nature* 429: 65-67.
 25. Bussiere LF, Hunt J, Jennions MD, Brooks R (2006) Sexual conflict and cryptic female choice in the black field cricket, *Teleogryllus commodus*. *Evolution* 60: 792-800.
 26. Bando T, Hamada Y, Kurita K, Nakamura T, Mito T, et al. (2011) Lowfat, a mammalian Lix1 homologue, regulates leg size and growth under the Dachous/Fat signaling pathway during tissue regeneration. *Developmental Dynamics* 240: 1440-1453.
 27. Lynch JA, Peel AD, Drechsler A, Averof M, Roth S (2010) EGF signaling and the origin of axial polarity among the insects. *Current Biology* 20: 1042-1047.
 28. Kainz F, Ewen-Campen B, Akam M, Extavour CG (2011) Delta/Notch signalling is not required for segment generation in the basally branching insect *Gryllus bimaculatus*. *Development* 138: 5015-5026.
 29. Mito T, Shinmyo Y, Kurita K, Nakamura T, Ohuchi H, et al. (2011) Ancestral functions of Delta/Notch signaling in the formation of body and leg segments in the cricket *Gryllus bimaculatus*. *Development* 138: 3823-3833.
 30. Mito T, Kobayashi C, Sarashina I, Zhang H, Shinahara W, et al. (2007) *even-skipped* has gap-like, pair-rule-like, and segmental functions in the cricket *Gryllus bimaculatus*, a basal, intermediate germ insect (Orthoptera). *Developmental Biology* 303: 202-213.
 31. Mito T, Inoue Y, Kimura S, Miyawaki K, Niwa N, et al. (2002) Involvement of *hedgehog*, *wingless*, and *dpp* in the initiation of proximodistal axis formation during the regeneration of insect legs, a verification of the modified boundary model. *Mechanisms of Development* 114: 27-35.

32. Nakamura T, Mito T, Bando T, Ohuchi H, Noji S (2007) Dissecting insect leg regeneration through RNA interference. *Cellular and Molecular Life Sciences* 65: 64-72.
33. Nakamura T, Mito T, Miyawaki K, Ohuchi H, Noji S (2008) EGFR signaling is required for re-establishing the proximodistal axis during distal leg regeneration in the cricket *Gryllus bimaculatus* nymph. *Developmental Biology* 319: 46-55.
34. Feuillet C, Leach JE, Rogers J, Schnable PS, Eversole K (2011) Crop genome sequencing: lessons and rationales. *Trends in Plant Science* 16: 77-88.
35. Gregory TR (2005) Synergy between sequence and size in large-scale genomics. *Nature Reviews Genetics* 6: 699-708.
36. Schatz MC, Delcher AL, Salzberg SL (2010) Assembly of large genomes using second-generation sequencing. *Genome Research* 20: 1165-1173.
37. Hanrahan SJ, Johnston JS (2011) New genome size estimates of 134 species of arthropods. *Chromosome Res* 19: 809-823.
38. Takagi A, Kurita K, Terasawa T, Nakamura T, Bando T, et al. (2012) Functional analysis of the role of *eyes absent* and *sine oculis* in the developing eye of the cricket *Gryllus bimaculatus*. *Development Growth and Differentiation* 54: 227-240.
39. Nakamura T, Yoshizaki M, Ogawa S, Okamoto H, Shinmyo Y, et al. (2010) Imaging of transgenic cricket embryos reveals cell movements consistent with a syncytial patterning mechanism. *Current Biology* 20: 1641-1647.
40. Meier T, Chabaud F, Reichert H (1991) Homologous patterns in the embryonic development of the peripheral nervous system in the grasshopper *Schistocerca gregaria* and the fly *Drosophila melanogaster*. *Development* 112: 241-253.
41. Meier T, Reichert H (1991) Serially homologous development of the peripheral nervous system in the mouthparts of the grasshopper. *Journal of Comparative Neurology* 305: 201-214.
42. Meier T, Reichert H (1989) Embryonic Development and Evolutionary Origin of the Orthopteran Auditory Organs. *Journal of Neurobiology* 21: 592-610.
43. Takahashi T, Hamada A, Miyawaki K, Matsumoto Y, Mito T, et al. (2009) Systemic RNA interference for the study of learning and memory in an insect. *Journal of Neuroscience Methods* 179: 9-15.
44. Miyawaki K, Mito T, Sarashina I, Zhang H, Shinmyo Y, et al. (2004) Involvement of Wingless/Armadillo signaling in the posterior sequential segmentation in the cricket, *Gryllus bimaculatus* (Orthoptera), as revealed by RNAi analysis. *Mechanisms of Development* 121: 119-130.
45. Watanabe T, Ochiai H, Sakuma T, HOrch HW, Hamaguchi N, et al. (2012) Non-transgenic genome modifications in a hemimetabolous insect using zinc-finger and TAL effector nucleases. *Nature Communications* in press.
46. Inoue Y, Niwa N, Mito T, Ohuchi H, Yoshioka H, et al. (2002) Expression patterns of *hedgehog*, *wingless*, and *decapentaplegic* during gut formation of *Gryllus bimaculatus* (cricket). *Mechanisms of Development* 110: 245-248.
47. Kumar S, Blaxter ML (2010) Comparing *de novo* assemblers for 454 transcriptome data. *BMC Genomics* 11: 571.
48. Zeng V, Villanueva KE, Ewen-Campen B, Alwes F, Browne WE, et al. (2011) *De novo* assembly and characterization of a maternal and developmental

- transcriptome for the emerging model crustacean *Parhyale hawaiiensis*. BMC Genomics 12: 581.
49. Kanehisa M, Goto S, Furumichi M, Tanabe M, Hirakawa M (2010) KEGG for representation and analysis of molecular networks involving diseases and drugs. Nucleic Acids Research 38: D355-360.
 50. Zeng V, Extavour CG (2012) ASGARD: an open-access database of annotated transcriptomes for emerging model arthropod species. Database 2012: bas048.
 51. Iseli C, Jongeneel CV, Bucher P (1999) ESTScan: a program for detecting, evaluating, and reconstructing potential coding regions in EST sequences. Proceedings of the International Conference on Intelligent Systems for Molecular Biology: 138-148.
 52. Zdobnov EM, Apweiler R (2001) InterProScan--an integration platform for the signature-recognition methods in InterPro. Bioinformatics 17: 847-848.
 53. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, et al. (2005) InterProScan: protein domains identifier. Nucleic Acids Research 33: W116-120.
 54. Bru C, Courcelle E, Carrere S, Beausse Y, Dalmar S, et al. (2005) The ProDom database of protein domain families: more emphasis on 3D. Nucleic Acids Res 33: D212-215.
 55. Attwood TK, Bradley P, Flower DR, Gaulton A, Maudling N, et al. (2003) PRINTS and its automatic supplement, prePRINTS. Nucleic Acids Res 31: 400-402.
 56. Letunic I, Goodstadt L, Dickens NJ, Doerks T, Schultz J, et al. (2002) Recent improvements to the SMART domain-based sequence annotation resource. Nucleic Acids Res 30: 242-244.
 57. Haft DH, Selengut JD, White O (2003) The TIGRFAMs database of protein families. Nucleic Acids Res 31: 371-373.
 58. Bateman A, Coin L, Durbin R, Finn RD, Hollich V, et al. (2004) The Pfam protein families database. Nucleic Acids Res 32: D138-141.
 59. Hulo N, Sigrist CJ, Le Saux V, Langendijk-Genevaux PS, Bordoli L, et al. (2004) Recent improvements to the PROSITE database. Nucleic Acids Res 32: D134-137.
 60. Wu CH, Nikolskaya A, Huang H, Yeh LS, Natale DA, et al. (2004) PIRSF: family classification system at the Protein Information Resource. Nucleic Acids Res 32: D112-114.
 61. Gough J, Karplus K, Hughey R, Chothia C (2001) Assignment of homology to genome sequences using a library of hidden Markov models that represent all proteins of known structure. J Mol Biol 313: 903-919.
 62. Pearl FM, Lee D, Bray JE, Sillitoe I, Todd AE, et al. (2000) Assigning genomic sequences to CATH. Nucleic Acids Research 28: 277-282.
 63. Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J, et al. (2005) The PANTHER database of protein families, subfamilies, functions and pathways. Nucleic Acids Research 33: D284-288.
 64. Bendtsen JD, Nielsen H, von Heijne G, Brunak S (2004) Improved prediction of signal peptides: SignalP 3.0. Journal of Molecular Biology 340: 783-795.
 65. Sonnhammer EL, von Heijne G, Krogh A (1998) A hidden Markov model for predicting transmembrane helices in protein sequences. Proceedings of the

- International Conference on Intelligent Systems for Molecular Biology 6: 175-182.
66. WELCH BL (1947) The generalisation of student's problems when several different population variances are involved. *Biometrika* 34: 28-35.
67. Roche (2011) 454 Sequencing System Software Manual version 2.6. Part C: GS De Novo Assembler, GS Reference Mapper, SFF Tools: Roche.
68. Ewing B, Hillier L, Wendl MC, Green P (1998) Base-calling of automated sequencer traces using phred. I. Accuracy assessment. *Genome Research* 8: 175-185.
69. Ewing B, Green P (1998) Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Research* 8: 186-194.
70. O'Neil ST, Dzurisin JDK, Carmichael RD, Lobo NF, Emrich SJ, et al. (2010) Population-level transcriptome sequencing of nonmodel organisms *Erynnis propertius* and *Papilio zelicaon*. *BMC Genomics* 11: 310.
71. Chen S, Yang P, Jiang F, Wei Y, Ma Z, et al. (2010) *De novo* analysis of transcriptome dynamics in the migratory locust during the development of phase traits. *PLoS ONE* 5: e15633.
72. Ma Z, Yu J, Kang L (2006) LocustDB: a relational database for the transcriptome and biology of the migratory locust (*Locusta migratoria*). *BMC Genomics* 7: 11.
73. Kang L, Chen X, Zhou Y, Liu B, Zheng W, et al. (2004) The analysis of large-scale gene expression correlated to the phase changes of the migratory locust. *Proc Natl Acad Sci U S A* 101: 17611-17615.
74. Danley PD, Mullen SP, Liu F, Nene V, Quackenbush J, et al. (2007) A cricket Gene Index: a genomic resource for studying neurobiology, speciation, and molecular evolution. *BMC Genomics* 8: 109.
75. Badisco L, Huybrechts J, Simonet G, Verlinden H, Marchal E, et al. (2011) Transcriptome analysis of the desert locust central nervous system: production and annotation of a *Schistocerca gregaria* EST database. *PLoS ONE* 6: e17274.
76. Lu HL, Tanguy S, Rispe C, Gauthier JP, Walsh T, et al. (2011) Expansion of genes encoding piRNA-associated Argonaute proteins in the pea aphid: diversification of expression profiles in different plastic morphs. *PLoS ONE* 6: e28051.
77. Lin G-w, Chang C-c (2009) Cloning and developmental characterization of four *vasa* genes, *Apvasal-4*, in the parthenogenetic and viviparous pea aphid *Acyrtosiphon pisum*. *Mechanisms of Development* 126: S252 215-P017.
78. Shigenobu S, Bickel RD, Brisson JA, Butts T, Chang C-C, et al. (2010) Comprehensive survey of developmental genes in the pea aphid, *Acyrtosiphon pisum*: frequent lineage-specific duplications and losses of developmental genes. *Insect Molecular Biology* 19 Suppl 2: 47-62.
79. Pfreundt U, James DP, Tweedie S, Wilson D, Teichmann SA, et al. (2010) FlyTF: improved annotation and enhanced functionality of the *Drosophila* transcription factor database. *Nucleic Acids Research* 38: D443-447.
80. Kopan R, Ilagan MX (2009) The canonical Notch signaling pathway: unfolding the activation mechanism. *Cell* 137: 216-233.
81. Kitisin K, Saha T, Blake T, Golestaneh N, Deng M, et al. (2007) Tgf-Beta signaling in development. *Science STKE* 2007: cm1.
82. Nusse R, Varmus H (2012) Three decades of Wnts: a personal perspective on how a scientific field developed. *EMBO Journal* 31: 2670-2684.

83. Schindler C, Levy DE, Decker T (2007) JAK-STAT signaling: from interferons to cytokines. *J Biol Chem* 282: 20059-20063.
84. Shaul YD, Seger R (2007) The MEK/ERK cascade: from signaling specificity to diverse functions. *Biochimica et Biophysica Acta* 1773: 1213-1226.
85. Ingham PW, Nakano Y, Seger C (2011) Mechanisms and functions of Hedgehog signalling across the metazoa. *Nature Reviews Genetics* 12: 393-406.
86. Tordjmann T (2011) Hippo signalling: Liver size regulation and beyond. *Clinics and Research in Hepatology and Gastroenterology*.
87. Liu D, Finley RL, Jr. (2010) Cyclin Y is a novel conserved cyclin essential for development in *Drosophila*. *Genetics* 184: 1025-1035.
88. Yandell M, Mungall CJ, Smith C, Prochnik S, Kaminker J, et al. (2006) Large-scale trends in the evolution of gene structures within 11 animal genomes. *PLoS Comput Biol* 2: e15.
89. Chen F, Mackey AJ, Vermunt JK, Roos DS (2007) Assessing performance of orthology detection strategies applied to eukaryotic genomes. *PLoS One* 2: e383.
90. Attwood TK, Coletta A, Muirhead G, Pavlopoulou A, Philippou PB, et al. (2012) The PRINTS database: a fine-grained protein sequence annotation and analysis resource--its status in 2012. *Database* 2012: bas019.
91. Jaroszewski L, Li Z, Krishna SS, Bakolitsa C, Wooley J, et al. (2009) Exploration of uncharted regions of the protein universe. *PLoS Biology* 7: e1000205.
92. Ma C, Yang P, Jiang F, Chapuis MP, Shali Y, et al. (2012) Mitochondrial genomes reveal the global phylogeography and dispersal routes of the migratory locust. *Molecular Ecology*.
93. Ceccato P, Cressman K, Giannini A, Trzaska S (2007) The Desert Locust Upsurge in West Africa (2003-2005): Information on The Desert Locust Early Warning System, and The Prospects for Seasonal Climate Forecasting. *International Journal of Pest Management* 53: 7-13.
94. Lomer CJ, Bateman RP, Johnson DL, Langewald J, Thomas M (2001) Biological control of locusts and grasshoppers. *Annual Review of Entomology* 46: 667-702.
95. Nolte DJ (1963) A pheromone for melanization of locusts. *Nature* 200: 660-661.
96. Shaw KL, Parsons YM, Lesnick SC (2007) QTL analysis of a rapidly evolving speciation phenotype in the Hawaiian cricket *Laupala*. *Molecular Ecology* 16: 2879-2892.
97. Ellison CK, Wiley C, Shaw KL (2011) The genetics of speciation: genes of small effect underlie sexual isolation in the Hawaiian cricket *Laupala*. *Journal of Evolutionary Biology* 24: 1110-1119.
98. Wheeler DL, Barrett T, Benson DA, Bryant SH, Canese K, et al. (2007) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res* 35: D5-12.
99. Büning J (1994) *The Insect Ovary: ultrastructure, previtellogenic growth and evolution*. London: Chapman and Hall. 400 p.
100. Kainz F (2009) *Cell communication during patterning: Notch and FGF signalling in Gryllus bimaculatus and their role in segmentation [PhD]*. Cambridge: University of Cambridge. 161 p.
101. Bateman A, Coin L, Durbin R, Finn RD, Hollich V, et al. (2004) The Pfam protein families database. *Nucleic Acids Research* 32: D138-141.

102. Letunic I, Goodstadt L, Dickens NJ, Doerks T, Schultz J, et al. (2002) Recent improvements to the SMART domain-based sequence annotation resource. *Nucleic Acids Research* 30: 242-244.

Figure Legends

Figure 1. Oogenesis and embryogenesis in the cricket model organism *Gryllus bimaculatus*. (A) Adult female cricket perched on a gloved human finger for perspective. (B) Anterior tip of a single ovariole from an adult female ovary, showing oocytes (o) at early previtellogenic stages of oogenesis. A single large germinal vesicle (gv) is distinguishable in each oocyte. Unlike meroistic (containing nurse cells) *Drosophila* ovaries, *G. bimaculatus* ovaries are panoistic and lack nurse cells [99]. (C) A single late stage oocyte with a single layer of columnar follicle cells (fc). (D-J) Chronological stages of *G. bimaculatus* embryogenesis showing the range of embryonic stages represented in the transcriptome presented here. (D) A fertilized egg just after laying. The egg nucleus is distinguishable as a dense patch in the dorsal yolk (arrowhead). Ages are shown as days (d) after egg-laying at 29°C. (E-I) are 3D reconstructions of confocal optical sections of Hoechst 3342-stained embryos dissected free from the egg; (J) is a micrograph of a live embryo dissected free from the chorion. Abbreviations: A = abdomen; C = cerci; E = eye; H = head; G = gnathal segments; L1 = first thoracic leg; L2 = second thoracic leg; L3 = third thoracic leg; T = thorax. Scale bar is 100 µm in (B, C, E-I) and 500 µm in (D, J). Anterior is to the left in all panels. Photo in (A) courtesy of David Behl; photos in (D) and (J) from [100].

Figure 2. Distribution of average coverage (bp/contig) within contigs produced by *de novo* assembly of the *G. bimaculatus* transcriptome. The coverage within contigs is

calculated by dividing the total number of base pairs contained in the reads used to construct a contig by the length of that contig.

Figure 3. Assessment of gene discovery and read length capacity of the *G.*

***bimaculatus de novo* assembled transcriptome.** (A) Randomly selected subsets of the trimmed reads were assembled using Newbler v2.5 in 10% increments, up to and including 100% of trimmed reads. For each subassembly, the number of unique BLAST hits against the NCBI non-redundant database (**nr**) with an E-value cutoff of 1e-10 (red; left axis) and the average coverage per base pair (blue; right axis) was calculated (see text for details). The number of unique BLAST hits did not increase after at least 90% of reads (3,795,085 reads) were assembled, while the coverage per base pair continued to increase as reads were added to the assembly. (B) Isotig length distribution for each subassembly created as described in (A). (C) Isotig length distribution of each subassembly for isotigs ≥ 4 kb. High numbers (≥ 50) of isotigs over 4kb in length are achieved only when $\geq 40\%$ of reads (1,686,646 reads) are assembled.

Figure 4. Ortholog hit ratio analysis of the *G. bimaculatus de novo* assembled

transcriptome. The ortholog hit ratio is a comparison of the length of an assembled sequence to the total length of the full length transcript of its putative ortholog [70]. Values close to one suggest that a transcript predicted by the *de novo* assembly is close to full length. Ortholog hit ratios for the *G. bimaculatus* transcriptome sequences are compared to those for the previously reported *de novo* assembled transcriptome of another insect, the milkweed bug *Oncopeltus fasciatus* [11]. (A) Ortholog hit ratio

analysis of assembled isotigs. A majority (63.8%) of all *G. bimaculatus* isotigs (black bars) have an ortholog hit ratio of ≥ 0.5 (blue arrowhead), and 40.0% have an ortholog hit ratio of ≥ 0.8 (red arrowhead). These values are higher than those obtained for the *O. fasciatus de novo* assembled transcriptome (grey bars) [11]. (B) Ortholog hit ratio analysis of unassembled singletons. As expected, singletons represent much smaller proportions of putative full-length transcripts. 6.3% of *G. bimaculatus* singletons (black) have an ortholog hit ratio of ≥ 0.5 (blue arrowhead), while 0.8% have an ortholog hit ratio of ≥ 0.8 (red arrowhead). As for the isotig analysis, these values are higher than those obtained for the *O. fasciatus de novo* assembled transcriptome (grey) [11].

Figure 5. Phylogenetic comparison of proportion of known proteomes represented in the *G. bimaculatus de novo* assembled transcriptome. The number (bold) and percentage (bold italics) of proteome sequences with a putative *G. bimaculatus* ortholog in the *de novo* transcriptome assembly is shown for selected animals with sequenced genomes (based on top BLAST hit, E-value cutoff $1e-5$). Proteomes were predicted from genome sequence sources as shown in Table S1. Numbers in large font in red and blue ovals indicate average proportion of sequences from all tested insect and deuterostome proteomes, respectively, represented in the *G. bimaculatus* transcriptome.

Figure 6. Sequence extension and gene discovery in the *G. bimaculatus* Hedgehog and Hippo pathways. (A) The *de novo* transcriptome assembly of *G. bimaculatus* newly identifies most members of the *hedgehog* pathway (red), from which only the *hedgehog* ligand (blue) was previously known (GenBank accession AB044709). (B) The

transcriptome also adds significant sequence data to the fragments of many genes in the Hippo signaling pathway that had been previously identified (green). Seven genes of the known pathway were not identified in the transcriptome (yellow, white), two of which lack any sequence data in GenBank (white). GenBank accessions for previously identified sequences are as follows: *discs overgrown* (*dco*): AB443442; *expanded* (*ex*): AB378099; *warts* (*wts*): AB300574; *cyclin E* (*cycE*): AB378067; *hippo* (*hpo*): AB378070; *inhibitor of apoptosis protein* (*diap1*): AB378071; *mob as tumor suppressor* (*mats*): AB378072; *yorkie* (*yki*): AB378076; *scaffold protein salvador* (*sav*): AB378074; *Merlin* (*Mer*): AB378073; *Kibra*: DC445461.

Figure 7. Automated annotation of the *G. bimaculatus de novo* transcriptome assembly using Gene Predictor. (A) Comparison of the proportion of non-redundant assembly sequences, isotigs and singletons that obtained a significant BLAST hit against **nr** (black bars), and those that were assigned a putative orthology by Gene Predictor (GP; white bars), based on the best reciprocal top BLAST hit with the *Drosophila melanogaster* proteome (see Table S1). (B) Comparison of the proportion of sequences with a significant BLAST hit in **nr** that also had a putative orthology assignment based on Gene Predictor (dark grey bars). All sequences assigned putative orthologs by Gene Predictor also had significant BLAST hits in **nr** (light grey bars).

Figure 8. Coding region analysis of *G. bimaculatus de novo* transcriptome assembly sequences without significant BLAST hits in **nr.** Assembly products that failed to obtain significant BLAST hits in **nr** (white) were examined for the presence of coding

regions (green) using EST Scan [51]. Assembly sequences thus predicted to contain coding regions were examined for the presence of known coding domains (yellow) using InterPro Scan [52,53]. Results are shown separately for isotigs (A), singletons (B) and all non-redundant assembly products (C). See also Table 3.

Figure 9. Comparison of sequences lacking significant BLAST hits to nr, with *Laupala kohalensis* and *Locusta migratoria* databases. (A-C) Assembly products that failed to obtain significant BLAST hits to **nr** (white) were examined for significant similarity (magenta) to transcripts from at least one of *L. migratoria* or *L. kohalensis* [71,72,73,74]. (A'-C') Assembly sequences thus identified were parsed into sequences with significant hits among only *L. kohalensis* sequences (red), only *L. migratoria* sequences (blue), or both (yellow). Results are shown separately for isotigs (A, A'), singletons (B, A') and all non-redundant assembly products (C, A').

Figure 10. Principal protein domain composition of *G. bimaculatus* transcriptome sequences with highest similarity to *Laupala kohalensis* or *Locusta migratoria* sequences. Relative proportions of the top 25 protein domains coded by *G. bimaculatus* transcriptome sequences with significant similarity to sequences from *L. kohalensis* (A), *L. migratoria* (B), or sequences from **nr** (C). Protein domain nomenclature from Pfam [101] as follows: AdoHcyase_NAD: PF00670; Ank: PF00023; ATP-gua_Ptrans/N: PF02807; BTB/POZ: PF00651; C2: PF00168; DUF (combined): n/a; EFG domains (combined): n/a; efhand/like: PF09279; F-box: PF00646; Glyco_hydro (combined): n/a; GTP_EFTU domains: PF00009; Laps: PF10169; LRR_1: PF00560; Metallophos:

PF00149; Myb_DNA-binding (combined): n/a; OS-D: PF03392; PARP: PF00644;
PGAMP: PF07644; Pkinase: PF00069; Ras: PF00071; Ribosomal (combined): n/a;
RRM_1: PF00076; RVT_1: PF00078; ubiquitin: PF00240; zinc finger (combined): n/a.

“Combined” indicates that multiple Pfam accessions are combined.

Table 1. Large-scale Orthopteran transcriptome resources to date.

	<i>Locusta migratoria</i>¹	<i>Laupala kohalensis</i>²	<i>Schistocerca gregaria</i>³	<i>Locusta migratoria</i>⁴	<i>Gryllus bimaculatus</i>⁵
Orthopteran Suborder	Caelifera	Ensifera	Caelifera	Caelifera	Ensifera
Superfamily	Acridoidea	Grylloidea	Acridoidea	Acridoidea	Grylloidea
Family	Acrididae	Gryllidae	Acrididae	Acrididae	Gryllidae
Sequencing Platform	Sanger	Sanger	Sanger	Illumina	454 Titanium
Tissue Source(s)	L5 ⁶	L5-L8 CNS	L3-L5 & adult CNS	Mainly L4	Ovaries & embryos
Normalized Library	No	Yes	Yes	No	No
# Raw Reads	76,012	14,502	nd	447,718,464	4,248,346
# Reads Used in Assembly	45,449	14,377	34,672	nd	4,216,721
# bp Used in Assembly	21,760,812	10,121,408	nd	nd	1,449,059,795
% Raw Reads Assembled	59.79%	99.14%	nd	nd	99.26%
# Contigs or Clusters	4,550	2,575	4,785	72,977	43,321
N50⁷ or Mean Contig Length (bp)	471	935	750	2,275	2,133
# Singletons or # Single ESTs	7,611	6,032	7,924	nd	120,805
% Singletons (of assembled reads)	16.75%	41.96%	22.85%	nd	2.86%
# Total Assembly Products	12,161	8,607	12,709	72,977	142,317
# Unigenes or # Unique BLAST Hits to nr	12,616	8,607	12,709	11,490	19,874

1. Data from [72,73].

2. Data from [74].

3. Data from [75].

4. Data from [71].

5. Data from this report.

6. L= larval stage. nd = data not reported in the relevant publication [71,72,73,75].

7. “N50” refers to isotig N50 from the *G. bimaculatus de novo* transcriptome assembly; mean contig length is shown for all other orthopteran transcriptome resources in this table.
8. # singletons are shown for the *G. bimaculatus de novo* transcriptome assembly; # single ESTs (not incorporated into contigs) are shown for all other orthopteran transcriptome resources in this table.
9. # unique BLAST hits against **nr** are shown for the *G. bimaculatus de novo* transcriptome assembly; # unigenes are shown for all other orthopteran transcriptome resources in this table.

Table 2. Assembly statistics and BLAST results against nr for the *G. bimaculatus de novo* transcriptome assembly.

Parameter ¹	Value
# bp Raw reads	1,483,726,666
Maximum raw read length	803
Minimum raw read length	13
Median raw read length	364
Maximum assembled read length	771
Minimum assembled read length	20
Median assembled read length	358
# Isogroups ² ("genes")	16,456
Mean # isotigs per isogroup	1.2
# Isotigs	21,512
Maximum isotig length	10,865
Minimum isotig length	57
Median isotig length	1,054.5
# Isotigs with BLAST hit against nr ³ , E-value cutoff e-10 (% of all isotigs)	11,135 (51.8%)
# Isotigs with BLAST hit against nr , E-value cutoff e-5 (% of all isotigs)	11,943 (55.5%)
Mean # contigs per isotig	1.7
# Singletons	120,805
Maximum singleton length	620
Minimum singleton length	50
Median singleton length	250.5
# Singletons with BLAST hit against nr , E-value cutoff e-10 (% of all singletons)	7,914 (6.6%)
# Singletons with BLAST hit against nr , E-value cutoff e-5 (% of all singletons)	10,815 (9.0%)
# Non-redundant assembly products (NRAP)	142,317
# NRAP with BLAST hit against nr , E-value cutoff e-10 (% of all NRAP)	19,049 (13.4%)
# NRAP with BLAST hit against nr , E-value cutoff e-5 (% of all NRAP)	22,758 (16.0%)
Total # BLAST hits ⁴ (nr)	22,758
Average coverage/bp	51.3

1. Values for number of raw reads, number and % of raw reads assembled (passed quality filters described in main text), number of contigs, isotig N50, % of singletons, total number of assembly products, and number of unique BLAST hits against **nr**, are shown in Table 1.
2. Because isogroups are collections of isotigs that are hypothesized to originate from the same gene, they do not comprise a single sequence and so cannot be mapped to **nr** using BLAST.
3. **nr** = NCBI non-redundant database.
4. For BLAST against **nr** the E-value cutoff was 1e-5. For breakdown of BLAST hits among different classes of assembly sequences, see Table 3.

Table 3. Length parameters of isotigs according to BLAST annotation and predicted protein-coding status.

BLAST hit¹/predicted protein coding status	Parameter	Value
Significant hit against nr^{2,3}	Maximum sequence length ⁴	10865
	Minimum sequence length	91
	Median sequence length	1669.50
	Average sequence length	1927.98
Significant hit against nr and contains predicted protein-coding region(s)	Maximum sequence length	10865
	Minimum sequence length	168
	Median sequence length	1730.5
	Average sequence length	1997.42
	Maximum predicted peptide length ⁵	2076
	Minimum predicted peptide length	11
	Median predicted peptide length	317.50
	Average predicted peptide length	386.82
No significant hit against nr	Maximum sequence length	6886
	Minimum sequence length	57
	Median sequence length	728.50
	Average sequence length	924.277
No significant hit against nr and contains predicted protein-coding region(s)	Maximum sequence length	6686
	Minimum sequence length	60
	Median sequence length	858.5
	Average sequence length	1130.16
	Maximum predicted peptide length	1710
	Minimum predicted peptide length	7
	Median predicted peptide length	144.5
	Average predicted peptide length	197.61
All NRI⁶ containing predicted protein-coding regions	Maximum sequence length	10865
	Minimum sequence length	60
	Median sequence length	1544.50
	Average sequence length	1837.57
	Maximum predicted peptide length	2076
	Minimum predicted peptide length	7
	Median predicted peptide length	282.50
	Average predicted peptide length	351.95
All NRI without predicted protein-coding regions	Maximum sequence length	6677
	Minimum sequence length	57
	Median sequence length	708.50
	Average sequence length	878.27
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences⁷	Maximum sequence length	5287
	Minimum sequence length	124
	Median sequence length	1093.50
	Average sequence length	1358.21
	Maximum predicted peptide length	1710
	Minimum predicted peptide length	25
	Median predicted peptide length	244.50
	Average predicted peptide length	320.84
No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences⁸	Maximum sequence length	6677
	Minimum sequence length	62
	Median sequence length	1004.50
	Average sequence length	1304.64

	Maximum predicted peptide length	1710
	Minimum predicted peptide length	16
	Median predicted peptide length	248.50
	Average predicted peptide length	315.37

1. BLAST E-value cutoff is e-5 for all hits reported in this table.
2. **nr** = NCBI non-redundant database.
3. Numbers of sequences in each category are shown in Figure 9.
4. Sequence lengths are reported in base pairs.
5. Predicted peptide lengths are reported in amino acids.
6. NRI = all non-redundant isotigs regardless of BLAST results against **nr**.
7. *Locusta migratoria* sequences used for comparison are from [72,73].
8. *Laupala kohalensis* sequences used for comparison are from [74].

Table 4. Length parameters of singletons according to BLAST annotation and predicted protein-coding status.

BLAST hit ¹ /predicted protein coding status	Parameter	Value
Significant hit against nr^{2,3}	Maximum sequence length ⁴	582
	Minimum sequence length	66
	Median sequence length	340.00
	Average sequence length	334.25
Significant hit against nr and contains predicted protein-coding region(s)	Maximum sequence length	574
	Minimum sequence length	68
	Median sequence length	343.5
	Average sequence length	337.54
	Maximum predicted peptide length ⁵	192
	Minimum predicted peptide length	8
	Median predicted peptide length	103.50
No significant hit against nr	Average predicted peptide length	103.28
	Maximum sequence length	620
	Minimum sequence length	50
	Median sequence length	243.50
No significant hit against nr and contains predicted protein-coding region(s)	Average sequence length	251.67
	Maximum sequence length	586
	Minimum sequence length	50
	Median sequence length	231.5
	Average sequence length	243.16
	Maximum predicted peptide length	189
	Minimum predicted peptide length	5
	Median predicted peptide length	60.50
All NRS containing predicted protein-coding region(s)	Average predicted peptide length	65.02
	Maximum sequence length	586
	Minimum sequence length	50
	Median sequence length	255.5
	Average sequence length	268.89
	Maximum predicted peptide length	192
	Minimum predicted peptide length	5
	Median predicted peptide length	71.5
All NRS without predicted protein-coding regions	Average predicted peptide length	75.45
	Maximum sequence length	620
	Minimum sequence length	50
	Median sequence length	249.50
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences⁷	Average sequence length	255.51
	Maximum sequence length	552
	Minimum sequence length	52
	Median sequence length	299
	Average sequence length	283.97
	Maximum predicted peptide length	176
	Minimum predicted peptide length	17
	Median predicted peptide length	74.50
No significant hit against nr and significant hit against <i>Laupala kohalensis</i>	Average predicted peptide length	75.08
	Maximum sequence length	597
	Minimum sequence length	52
	Median sequence length	286.50

sequences ⁸	Average sequence length	280.55
	Maximum predicted peptide length	188
	Minimum predicted peptide length	11
	Median predicted peptide length	77.5
	Average predicted peptide length	77.40

1. BLAST E-value cutoff is e-5 for all hits reported in this table.
2. **nr** = NCBI non-redundant database.
3. Numbers of sequences in each category are shown in Figure 9.
4. Sequence lengths are reported in base pairs.
5. Predicted peptide lengths are reported in amino acids.
6. NRS = all non-redundant singletons regardless of BLAST results against **nr**.
7. *Locusta migratoria* sequences used for comparison are from [72,73].
8. *Laupala kohalensis* sequences used for comparison are from [74].

Table 5. Statistical comparison of isotig and singleton nucleotide sequence lengths according to BLAST annotation and predicted protein-coding status. Values shown are $p \geq 0.05$ value results of a Welch's t-test. *** = $p < 0.0001$; * $p < 0.05$.

BLAST hit ¹ /predicted protein coding status ²	Significant hit against nr ²	Significant hit against nr and contains predicted coding regions	No significant hit against nr	No significant hit against nr and contains predicted coding regions	All NRAS ³ containing predicted protein-coding regions	All NRAS without predicted protein-coding regions	No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences	No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences
ISOTIGS⁴								
Significant hit against nr		0.9998	***	***	***	***	***	***
Significant hit against nr and contains predicted coding regions			1	***	1	1	1	1
No significant hit against nr				1	1	***	1	1
No significant hit against nr and contains predicted coding regions					***	1	***	***
All NRAS containing predicted protein-coding regions						***	***	***
All NRAS without predicted protein-coding regions							1	1
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences								0.2268
No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences								
SINGLETONS								
Significant hit against nr		0.9798	***	***	***	***	***	***

Significant hit against nr and contains predicted coding regions			1	***	1	1	1	1
No significant hit against nr				***	1	***	0.4208	1
No significant hit against nr and contains predicted coding regions					***	***	*	***
All NRAS containing predicted protein-coding regions						***	***	0.0969
All NRAS without predicted protein-coding regions							0.1358	0.9985
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences								0.9967
Significant hit against <i>Laupala kohalensis</i> sequences								

1. BLAST E-value cutoff is e-5 for all hits reported in this table.
2. **nr** = NCBI non-redundant database.
3. NRAS = all non-redundant assembly products (isotigs or singletons) regardless of BLAST results against **nr**.
4. Numbers of sequences in each category are shown in Figure 9. Mean, median, maximum and minimum values for each category are shown in Tables 3 and 4.

Table 6. Statistical comparison of isotig and singleton predicted coding sequence lengths according to BLAST annotation status. Values shown are $p \geq 0.05$ value results of a Welch's t-test. *** = $p < 0.0001$; ** $p < 0.001$; * $p < 0.05$

BLAST hit ¹ /predicted protein coding status ²	Significant hit against nr ²	No significant hit against nr	All NRAS ³	No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences	No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences
ISOTIGS⁴					
Significant hit against nr		***	1	**	***
No significant hit against nr			***	1	1
All NRAS				*	0.0059
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences					0.4052
No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences					
SINGLETONS					
Significant hit against nr		***	1	***	***
No significant hit against nr			***	1	1
All NRAS				0.4091	0.9235
No significant hit against nr and significant hit against <i>Locusta migratoria</i> sequences					0.8685
No significant hit against nr and significant hit against <i>Laupala kohalensis</i> sequences					

1. BLAST E-value cutoff is e-5 for all hits reported in this table.
2. **nr** = NCBI non-redundant database.
3. NRAS = all non-redundant assembly products regardless of BLAST results against **nr**.
4. Numbers of sequences in each category are shown in Figure 9. Mean, median, maximum and minimum values for each category are shown in Tables 3 and 4.

Supporting Information: Supplementary Figure and Table Legends

Figure S1. Comparison of read lengths from *de novo* assembly of the *G. bimaculatus* transcriptome. (A) Distribution of read lengths before (black) and after (blue) trimming to remove low quality reads (see text for details). (B) Distribution of trimmed read lengths before (blue) and after (red) assembly with Newbler v2.5. The assembly yielded assembled reads of over 10,000 bp. (C) Distribution of read lengths of the shortest assembled (red) and raw (blue) reads.

Figure S2. Schematics of conserved metazoan signal transduction pathways showing components identified in the *G. bimaculatus* transcriptome. BLAST was used to search for signaling pathway genes in the *G. bimaculatus* transcriptome (see Table S4); genes with newly identified putative orthologs are indicated in red. Genes outlined in grey with grey typeface indicate genes without *D. melanogaster* homologs. Pathway schematics are modified from KEGG pathway model images (<http://www.genome.jp/kegg/kegg1.html>). (A) Notch pathway. (B) TGF β pathway. (C) Wnt pathway. (D) Janus Kinase (JAK)-signal transducer and activator of transcription (STAT) pathway. (E) Mitogen-activated protein Kinase (MAPK) pathway.

Figure S3. Complete protein domain composition of *G. bimaculatus* transcriptome sequences with highest similarity to *Laupala kohalensis* or *Locusta migratoria* sequences. Relative proportions of all protein domains coded by *G. bimaculatus* transcriptome sequences with significant similarity to sequences from *L. kohalensis* (A), *L. migratoria* (B), or sequences from **nr** (C). Protein domain nomenclature from Pfam

[101] and SMART [102] databases as follows: 5_nucleotid_C: PF2872; Abhydrolase_1: PF00561; adh_short: PF00106; ADK: OF00406; AdoHcyase_NAD: PF00670; Amidohydro_1: PF01979; Ank: PF00023; AP_endonuc_2_N: PF07582; Asparaginase_2: PF01112; ATP-gua_Ptrans/N: PF02807; BAH: PF01426; BTB/POZ: PF00651; Btz: SM 01044; bZIP_2: PF07716; C2: PF00168; CBM_14: PF01607; COesterase: PF00135; Cyclin_N: PF00134; Cys_Met_Meta_PP: PF01053; DEAD: PF00270; DUF (combined): n/a; EFG domains (combined): n/a; efhand/like: PF09279; eIF-5_eIF-2B: PF01873; ELM2: PF01448; ELO: PF01151; EMP70: PF02990; ETF_alpha: PF00766; Exo_endo_phos: PF03372; F-box: PF00646; fn3: PF00041; G-patch: PF01858; GATA: PF00320; GCV_H: PF01597; GHMP_kinases_N: PF00288; Glyco_hydro (combined): n/a; GTP_EFTU domains: PF00009; HECT: PF00632; Hemocyanin_N: PF03722; HSP90: PF00183; IF-2B: PF01008; IPP-2: PF04979; JHBP: PF06585; Laps: PF10169; Ldl_recept_a: PF00057; Lectin_C: PF00059; LRR_1: PF00560; MA3: PF00560; MADF_DNA_bdg: PF10545; MAP65_ASE1: PF03999; Metallophos: PF00149; MIF4G: PF02854; Myb_DNA-binding (combined): n/a; NAC: PF01849; NAP: PF00956; NDUF_B8: PF05821; NIPSNAP: PF07978; Nucleoplasmin: PF03066; OS-D: PF03392; p450: PF00067; PABP: PF00658; PARP: PF00644; Peptidase_M17: PF00883; PGAMP: PF07644; PH: PF00169; PI-PLC-X/Y: PF00378/8; Pkinase: PF00069; PTPS: PF01242; Ras: PF00071; Ribophorin_I: PF04597; Ribosomal (combined): n/a; RNA_pol_A_bac: PF01000; RnaseH: PF00075; RRM_1: PF00076; RVT_1: PF00078; SAM_1: PF00536; Sedlin_N: PF04628; Serpin: PF00079; SH2: PF00017; SH3_1: PF00018; SNase: PF00565; Stathmin: PF008310; Synaptobrevin: PF00957; Thioredoxin: PF00085; Thymosin: PF01290; TRAP-gamma: PF07074; TRM: PF02005; TUDOR: PF00567;

ubiquitin: PF00240; W2: PF02020; WD40: PF00400; zinc finger (combined): n/a.

“Combined” indicates that multiple Pfam accessions are combined.

Table S1. Sources of proteome sequences from animals with sequenced genomes used for comparison with the *G. bimaculatus de novo* transcriptome assembly.

Sequences were used for ortholog hit ratio analyses (Figure 3) and phylogenetic comparisons of proportion of proteome sequences for which putative *G. bimaculatus* orthologs were found (Figure 4).

Table S2. Contribution of the *G. bimaculatus* transcriptome to GenBank accessions.

Sequences of *G. bimaculatus* developmental genes from GenBank were used as a query to BLAST the *de novo* transcriptome assembly. Matches in the transcriptome were found among both assembled reads and singletons.

Table S3. FlyTF transcription factor orthologs identified in the *G. bimaculatus* transcriptome. BLAST (E-cutoff 1e-5) was used to search the *G. bimaculatus* transcriptome for orthologs to the transcription factors belonging to the FlyTF database [79].

Table S4. Selected signaling pathway genes identified in the *G. bimaculatus* transcriptome. Hit ID indicates if gene hits were found assembled reads (A) or singletons (S). Length (range) indicates the shortest and longest A or S hit sequences for each gene. Query organism was *D. melanogaster* for all cases.

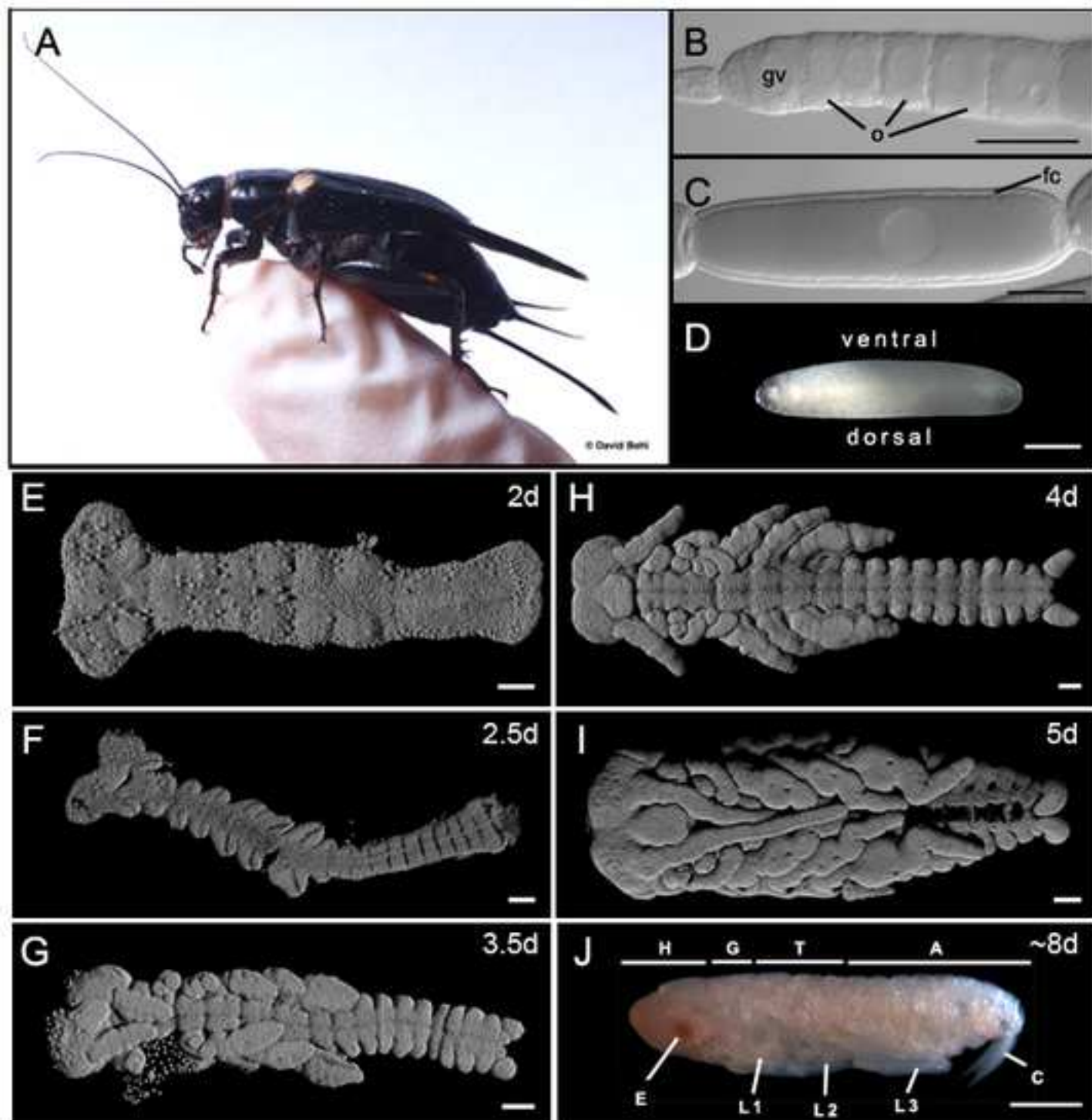
Table S5. Selected gametogenesis genes identified in the *G. bimaculatus*

transcriptome. Hit ID indicates if gene hits found were assembled reads (A) or singletons (S). Length (range) indicates the shortest and longest A or S hit sequences for each gene. Groups of hits of a given color indicate transcriptome sequences that mapped to the same overlapping region of the BLAST target (putative SNPs or isoforms); hits of different colors indicate transcriptome sequences that map to different, non-overlapping regions of the BLAST target. Query organism was *D. melanogaster* for all cases.

Table S6. Selected developmental process genes identified in the *G. bimaculatus de*

novo transcriptome assembly. Hit ID indicates if gene hits found were assembled reads (A) or singletons (S). Length (range) indicates the shortest and longest A or S hit sequences for each gene. Groups of hits of a given color indicate transcriptome sequences that mapped to the same overlapping region of the BLAST target (putative SNPs or isoforms); hits of different colors indicate transcriptome sequences that map to different, non-overlapping regions of the BLAST target. Query organism was *D. melanogaster* for all cases.

Zeng et al. Figure 1



Zeng et al. Figure 2

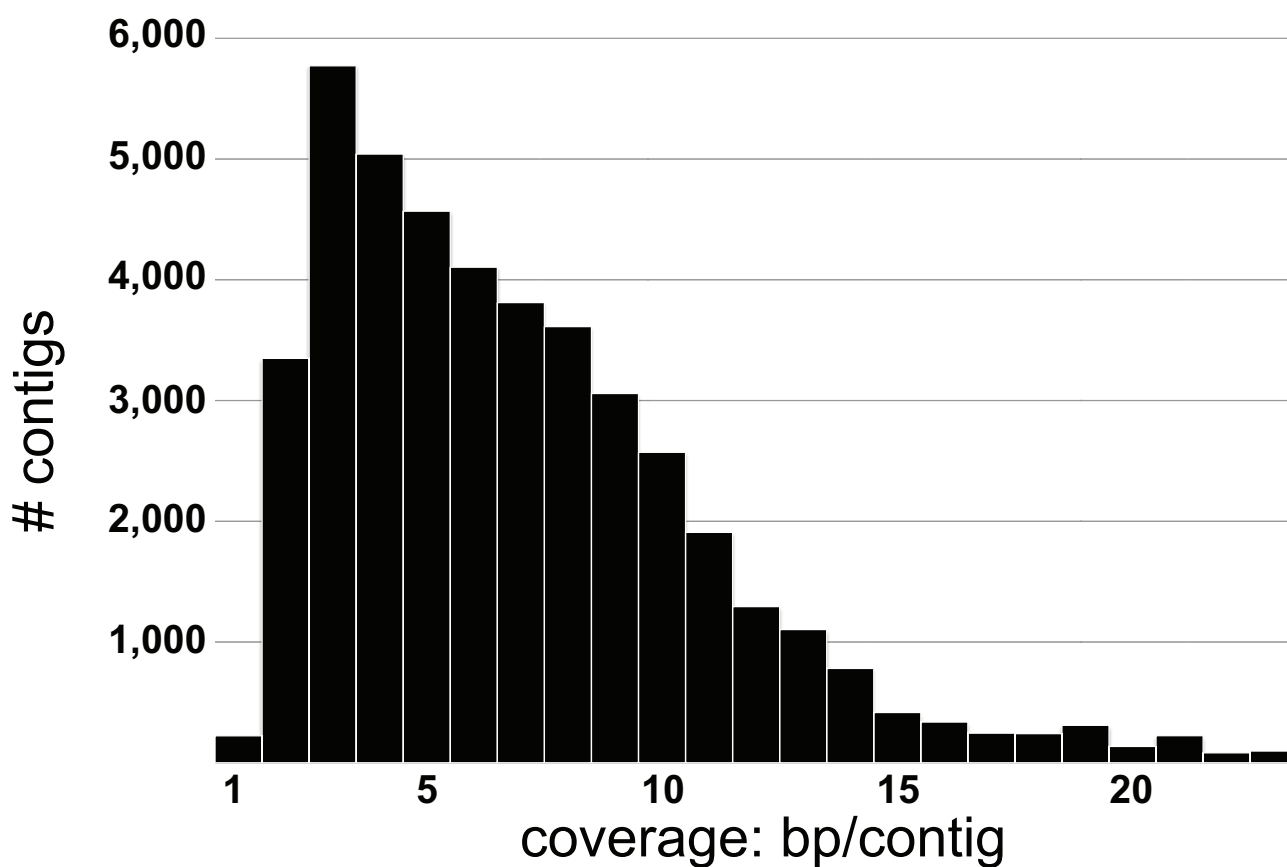


Figure 3

Zeng et al. Figure 3

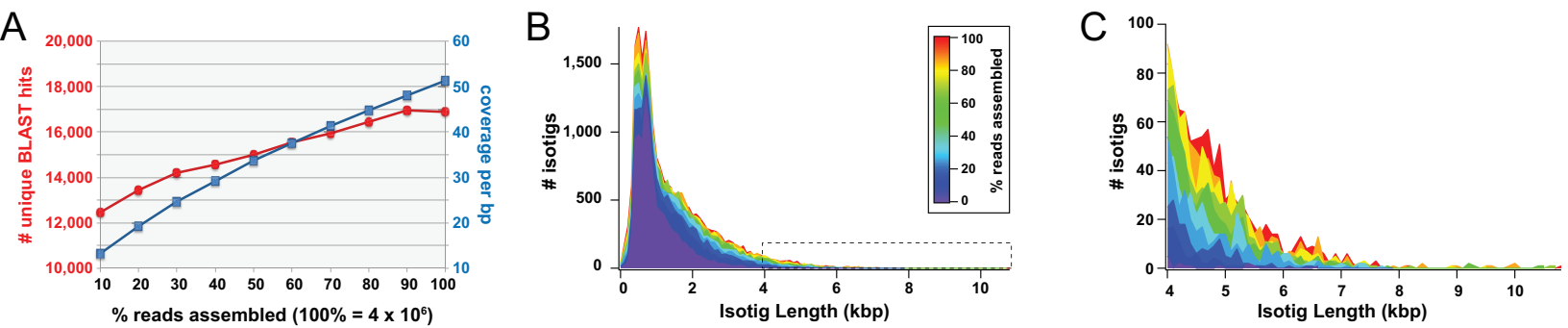
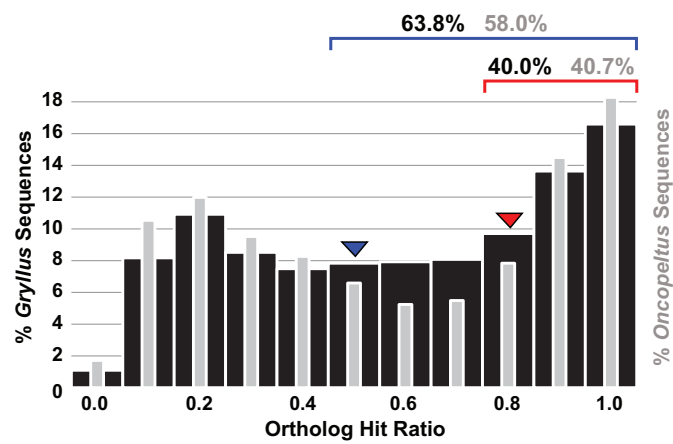


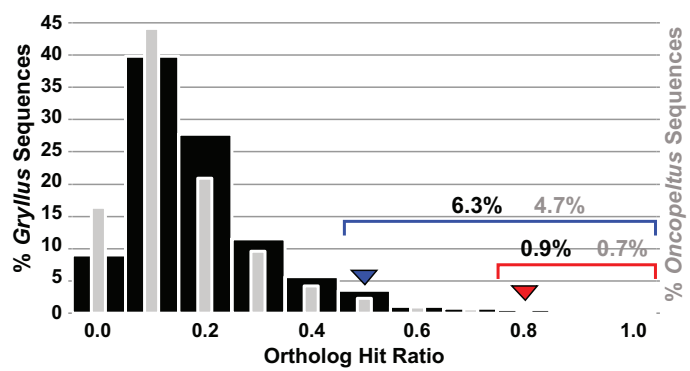
Figure 4

Zeng et al. Figure 4

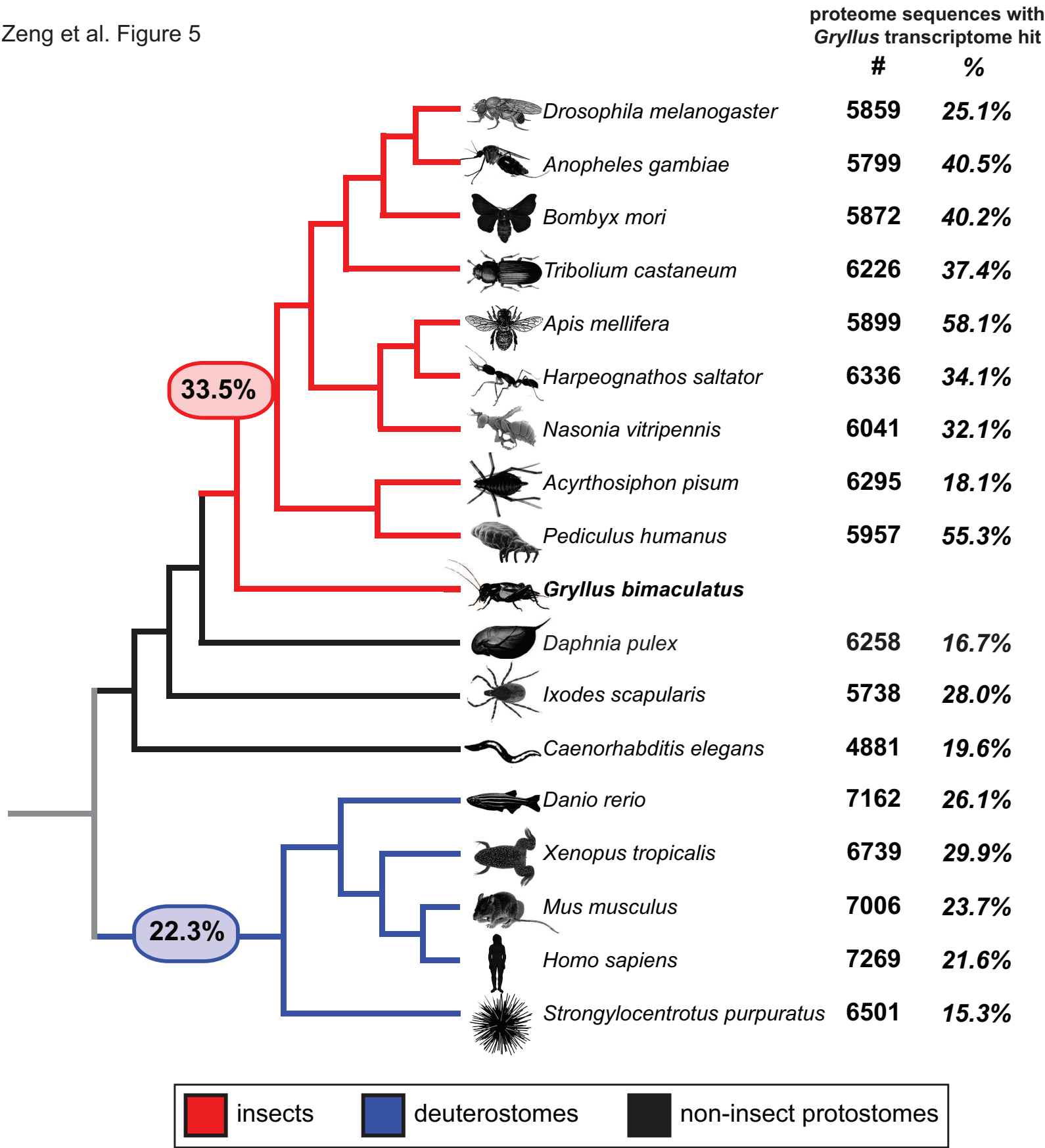
A isotigs



B singletons



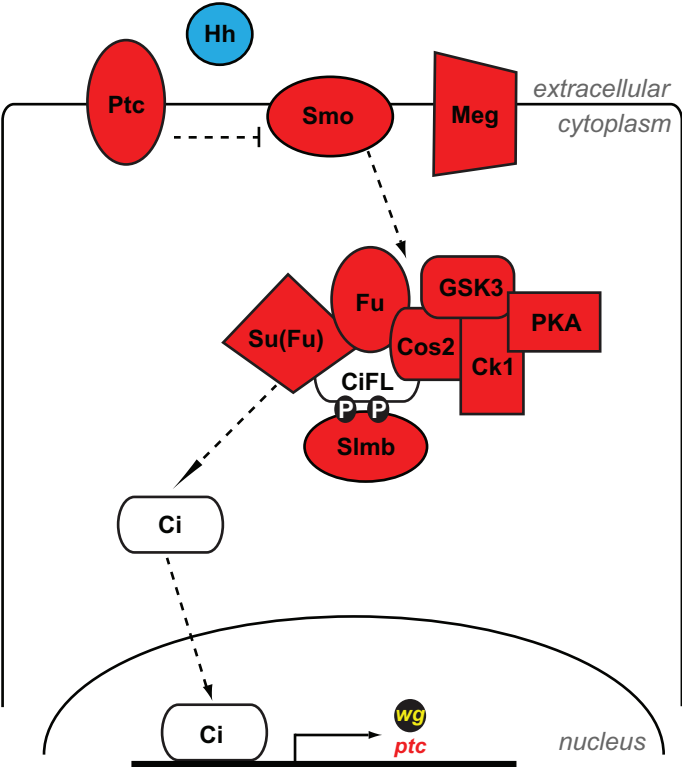
Zeng et al. Figure 5



Zeng et al. Figure 6

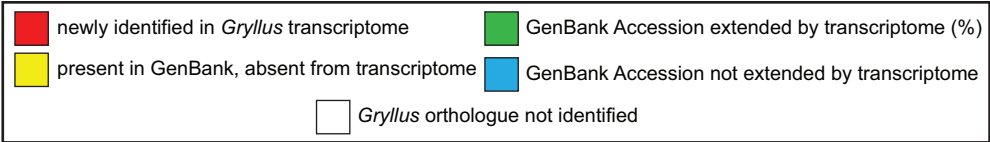
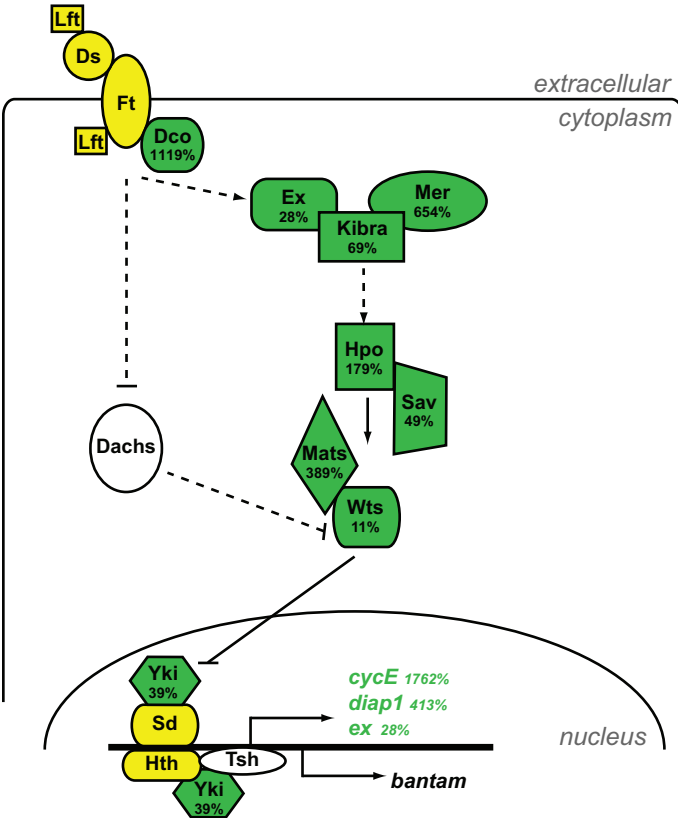
A

Hedgehog Signaling

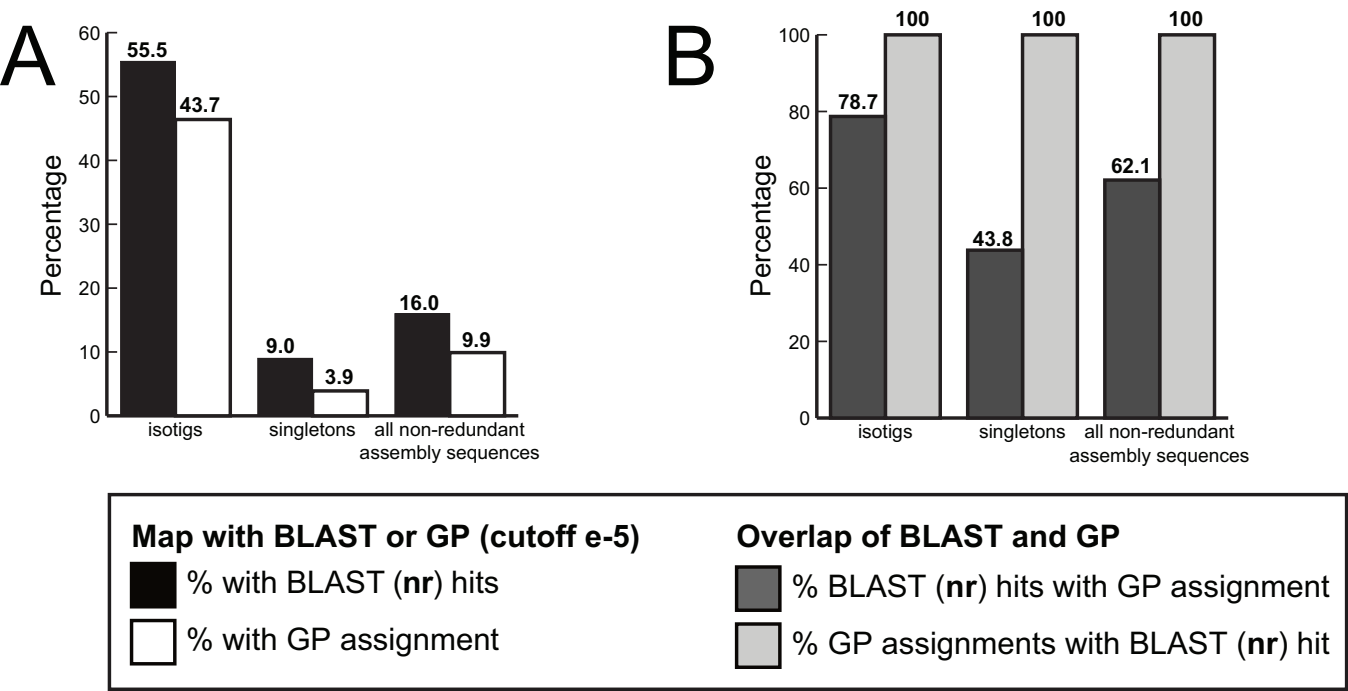


B

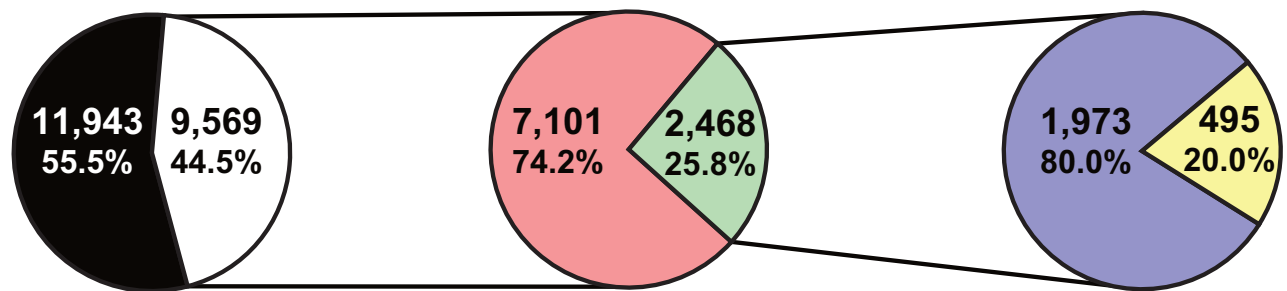
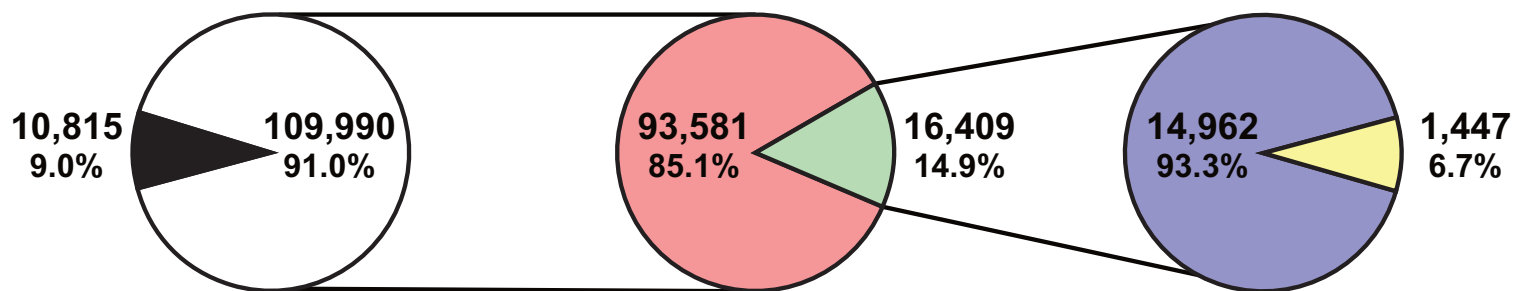
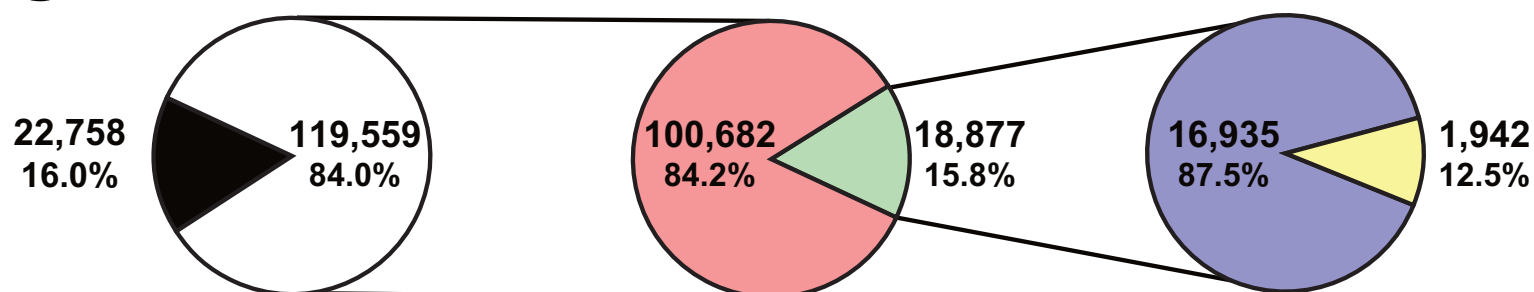
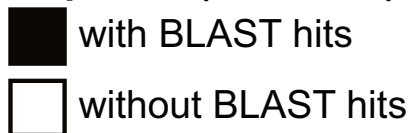
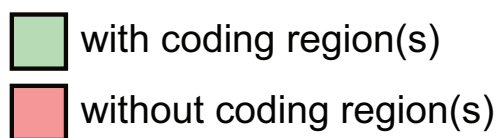
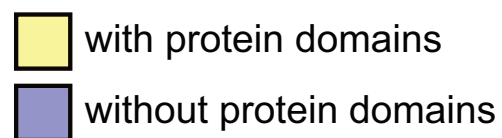
Hippo Signaling



Zeng et al. Figure 7



Zeng et al. Figure 8

A isotigs**B** singletons**C** all non-redundant assembly sequences**Map to nr (cutoff e-5)****EST Scan****InterPro Scan**

Zeng et al. Figure 9

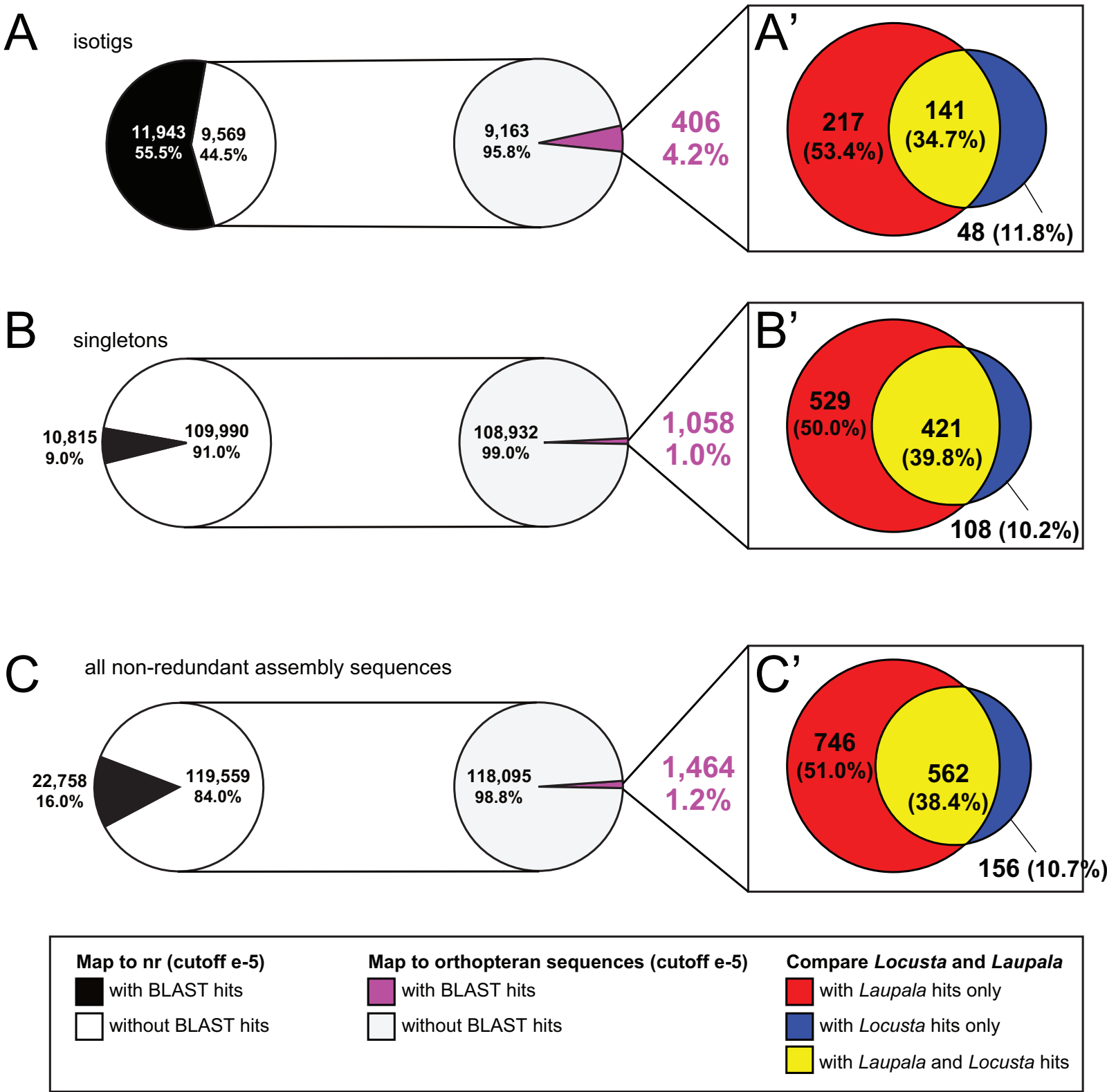
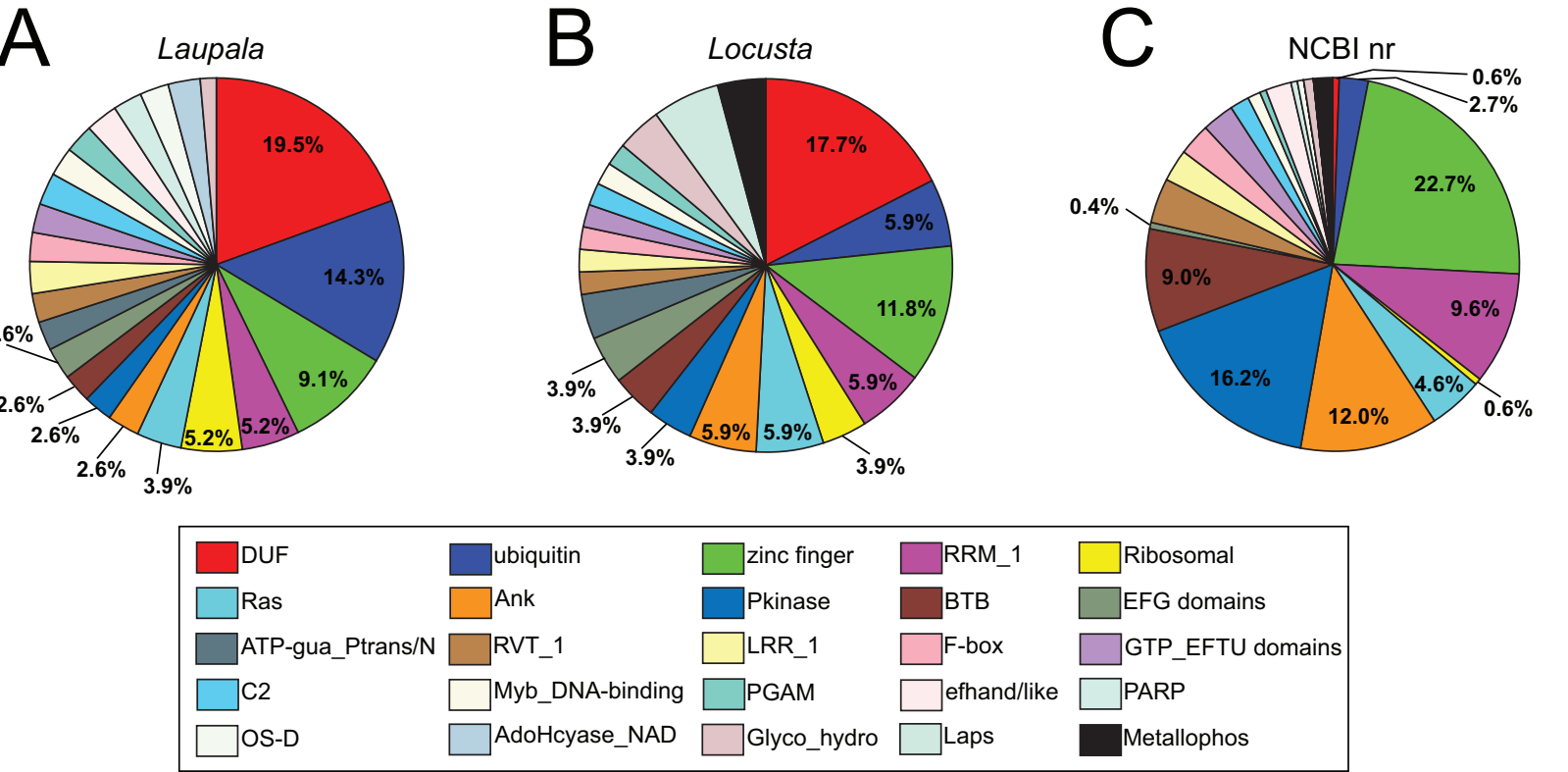
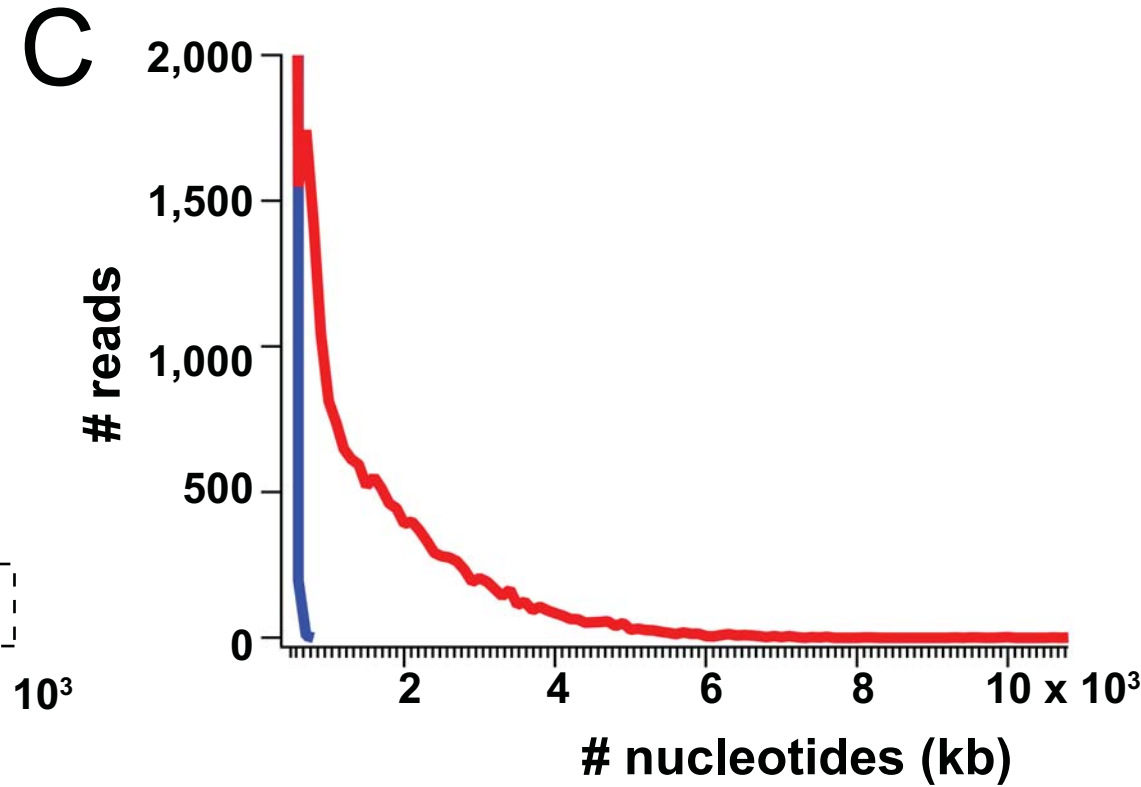
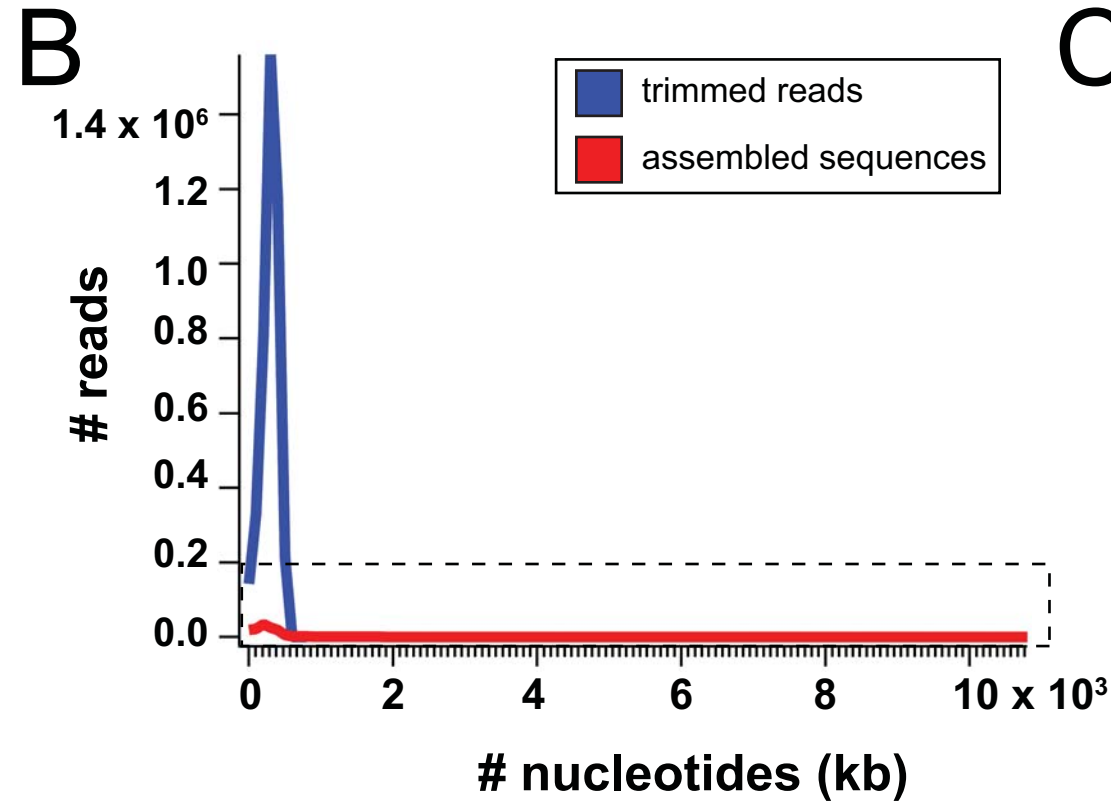
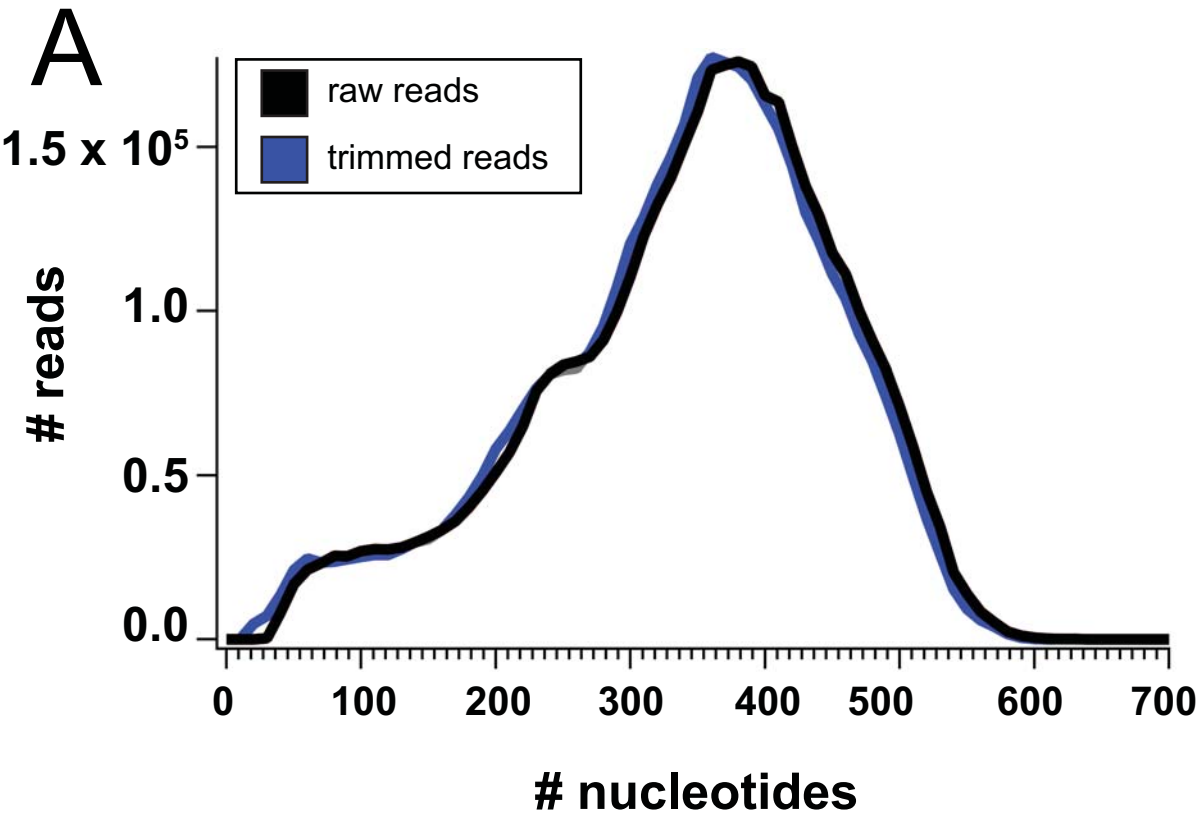


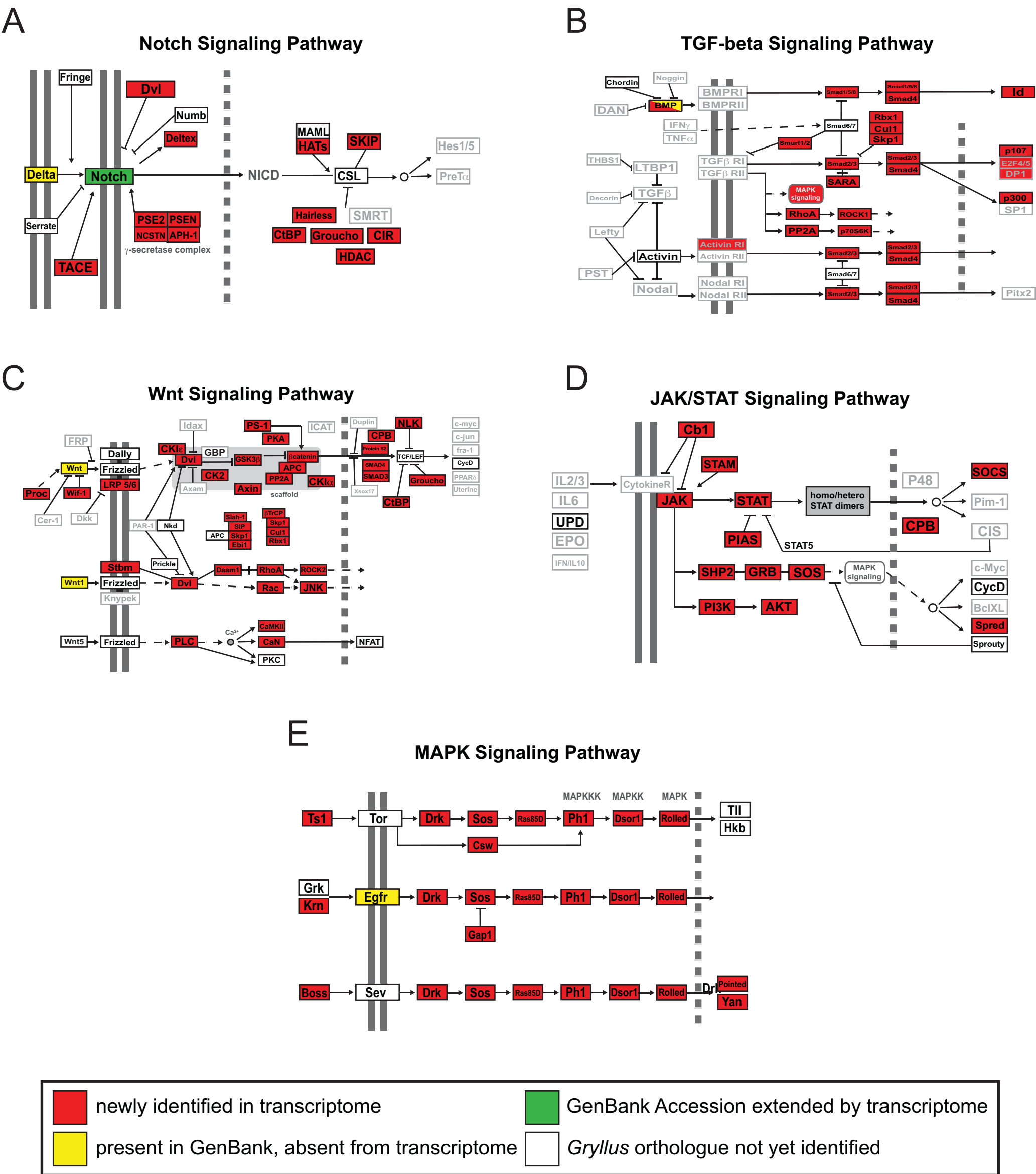
Figure 10

Zeng et al. Figure 10



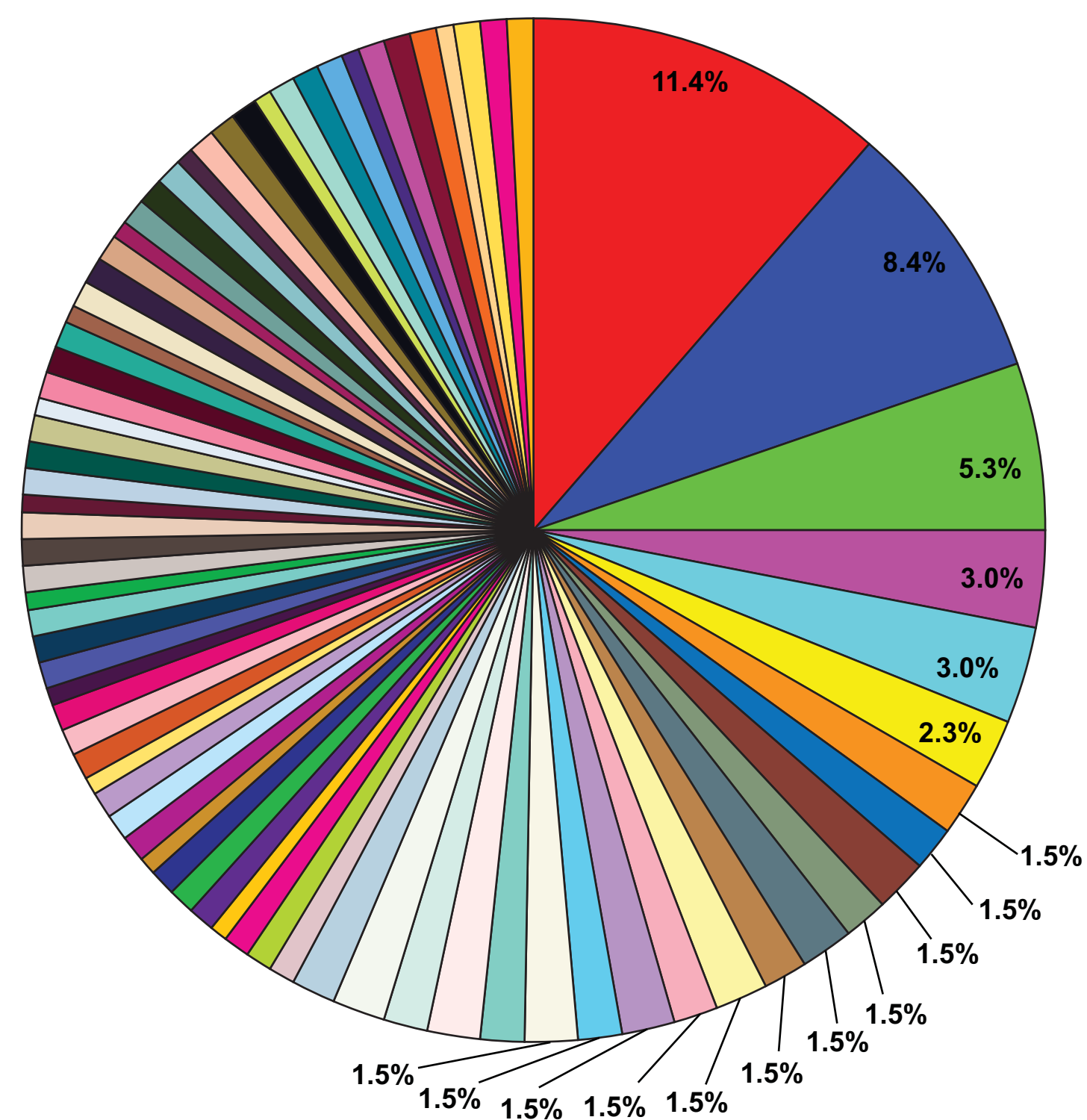


Zeng et al. Figure S2



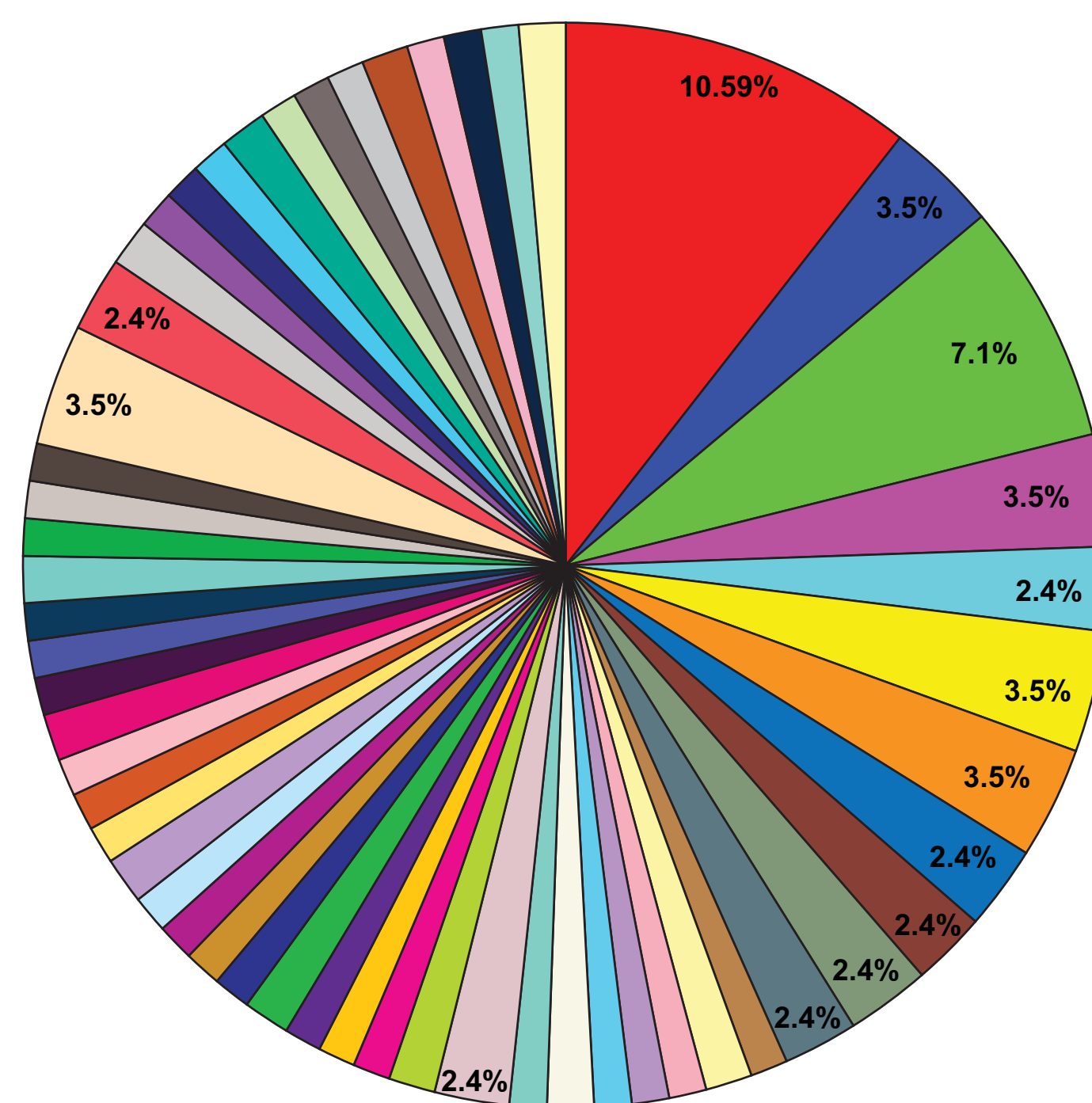
A

Laupala



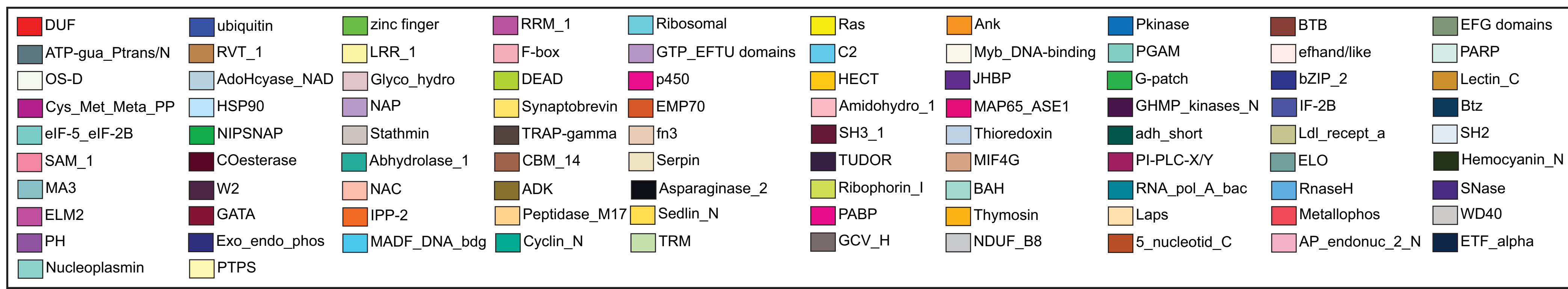
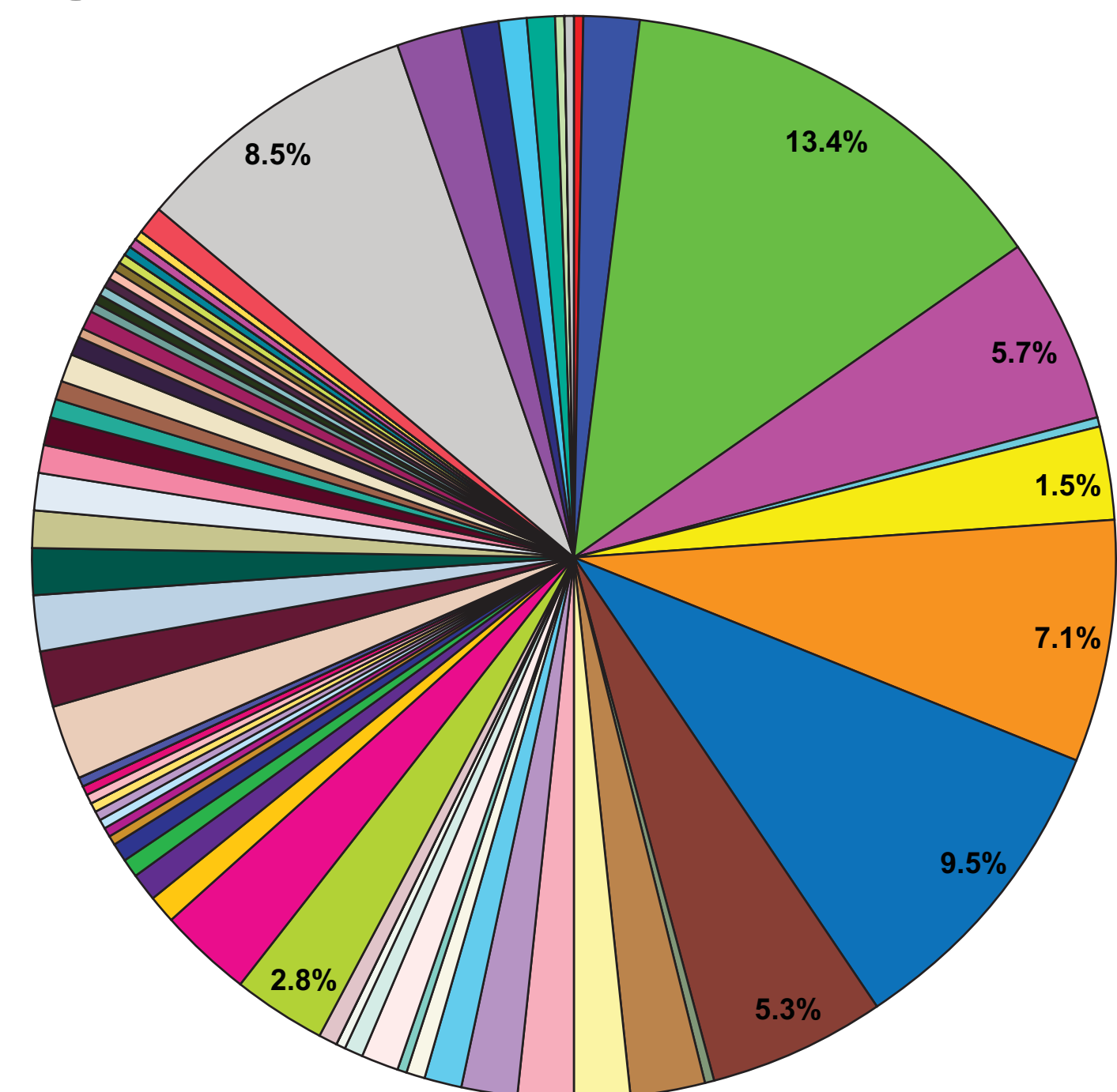
B

Locusta



C

NCBI nr



Sources of animal proteome data used for phylogenetic comparisons of *G. bimaculatus* transcriptome sequence matches

Species	Proteome Source	Download Date
<i>Apis mellifera</i>	http://hymenopteragenome.org/drupal/sites/hymenopteragenome.org.beebase/files/data/Amel_release1_OGS_pep.fa.gz	25Mar11
<i>Pediculus humanus</i>	ftp://ftp.vectorbase.org/public_data/organism_data/phumanus/Geneset/pediculus_humanus_PhumU1.2.fa.gz	25Mar11
<i>Anopheles gambiae</i>	ftp://ftp.vectorbase.org/public_data/organism_data/agambiae/Geneset/anopheles_gambiae_AgamP3.6.fa.gz	25Mar11
<i>Bombyx mori</i>	ftp://silkbdb.org/pub/current/Gene/silkpep.fa.gz	25Mar11
<i>Laupala kohalensis</i> ESTs	http://combio.dfci.harvard.edu	4May11
<i>Locusta migratoria</i> ESTs	http://locustdb.genomics.org.cn/download/Locust_EST.zip	4May11
<i>Tribolium castaneum</i>	ftp://bioinformatics.ksu.edu/pub/BeetleBase/3.0/Sequences/Tribolium_Official_Gene_Sequences/peptide.fa	25Mar11
<i>Camponotus floridanus</i>	http://hymenopteragenome.org/drupal/sites/hymenopteragenome.org.camponotus/files/data/cflo_v3.3.fa	25Mar11
<i>Saccharomyces cerevisiae</i>	http://downloads.yeastgenome.org/sequence/S288C_reference/orf_protein/orf_trans_all.fasta.gz	4May11
<i>Aedes aegypti</i>	ftp://ftp.vectorbase.org/public_data/organism_data/aaegypti/Geneset/aedes_aegypti_AaegL1.2.fa.gz	25Mar11
<i>Harpegnathos saltator</i>	http://hymenopteragenome.org/drupal/sites/hymenopteragenome.org.harpegnathos/files/data/hsal_v3.3.fa.gz	4May11
<i>Culex quinquefasciatus</i>	ftp://ftp.vectorbase.org/public_data/organism_data/cquinquefasciatus/Geneset/culex_quinquefasciatus_CpipJ1.2.fa.gz	25Mar11
<i>Gallus gallus</i>	ftp://ftp.ncbi.nih.gov/genomes/Gallus_gallus/protein/Gnomon_prot.fsa.gz	4May11
<i>Nasonia vitripennis</i>	http://genomes.arc.georgetown.edu/nasonia/nasonia_genome_consortium/data/Nvit_OGSv1.2_pep.fa.gz	25Mar11
<i>Xenopus tropicalis</i>	ftp://ftp.ncbi.nih.gov/refseq/X_tropicalis/mRNA_Prot/frog.protein.faa.gz	4May11
<i>Ixodes scapularis</i>	ftp://ftp.vectorbase.org/public_data/organism_data/iscapularis/Geneset/ixodes_scapularis_IscaW1.1.fa.gz	4May11
<i>Danio rerio</i>	ftp://ftp.ncbi.nih.gov/genomes/D_rerio/protein/Gnomon_prot.fsa.gz	4May11
<i>Drosophila melanogaster</i>	ftp://ftp.flybase.net/genomes/Drosophila_melanogaster/current/fasta/dmel-all-translation-r5.35.fasta.gz	25Mar11
<i>Mus musculus</i>	ftp://ftp.ncbi.nih.gov/genomes/M_musculus/protein/Gnomon_prot.fsa.gz	4May11
<i>Homo sapiens</i>	ftp://ftp.ncbi.nih.gov/refseq/H_sapiens/mRNA_Prot/human.protein.faa.gz	4May11
<i>Caenorhabditis elegans</i>	ftp://ftp.wormbase.org/pub/wormbase/species/c_elegans/sequence/protein/c_elegans.current.protein.fa.gz	4May11
<i>Acyrtosiphon pisum</i>	http://arthropods.eugenius.org/aphid/data/geneset1/ACYPIprot.fa.gz	25Mar11
<i>Daphnia pulex</i>	ftp://iubio.bio.indiana.edu/daphnia/genome/Daphnia_pulex/current/fasta/dpulex-all-translation-jgi060905.fasta.gz	11Dec11
<i>Strongylocentrotus purpuratus</i>	ftp://ftp.ncbi.nih.gov/genomes/Strongylocentrotus_purpuratus/protein/Gnomon_prot.fsa.gz	4May11
<i>Escherichia coli</i>	ftp://ftp.ncbi.nih.gov/genomes/Bacteria/Escherichia_coli_K_12_substr_DH10B_uid58979/NC_010473.faa	4May11

Table S2

Contribution of the *G. bimaculatus* transcriptome to existing *G. bimaculatus* GenBank accessions.

Gene Name	Accession #	Accession Sequence Length (nt)	total # nt added by transcriptome	# 5' nt added by transcriptome	# 3' nt added by transcriptome	% Accession lengthened by transcriptome	Transcriptome Read Name	Consensus Region	Query Location
14-3-3epsilon	AB443441	460	25	25	0	5%	GE8SX9M02IK8UO	1-365	1-367
		460	111	0	111	24%	GFCP6CO02GWC4H	425-112	151-460
14-3-3zeta	AB443440	438	2605	19	2586	595%	isotig03712	2597-3034	full
		438	2605	19	2586	595%	isotig03711	2597-2962	73-438
16s ribosomal	AF248685	498	207	207	0	42%	GE8SX9M02GQ4TQ	208-423, 493-533	1-215, 280-331
		498	120	0	120	24%	GE8SX9M01ELV4A	121-329	289-498
18S ribosomal	AF514548	1021	175	175	0	17%	isotig14176	176-799	1-627
		1021	938	0	938	92%	contig11156	939-1037	922-1021
28s ribosomal	EU878290	726	90	90	0	12%	isotig07604	91-377	1-287
		726	283	283	0	39%	isotig07603	284-570	1-287
		726	64	0	64	9%	isotig20138	65-144	645-724
abdominal-A	AB194277	868	309	309	0	36%	GFCP6CO01AM666	310-375, 386-458	708-780, 790-854
accessory gland protein actin	DQ630916	570	0	0	0	0%	GE8SX9M01BFM7Q	full	19-246
	AB087882	1290	0	0	0	0%	GFJY65E02JJW4Q	3-391	453-522, 874-1024, 1106-1275
aristaless	AB071147	2857	0	0	0	0%	GFCP6CO01BE1VK	full	2587-2831
armadillo protein beta-actin	AB109212	3836	160	160	0	4%	isotig05341	161-2579, 2604-3966	11-2428, 2453-3811
	DQ630919	210	203	138	65	97%	GFJY65E02ICT38	66-274	1-210
chico	AB370294	440	37	37	0	8%	GFJY65E01AQ60L	292-36	7-262
		440	0	0	0	0%	GFJY65E01CJM64	1-231	35-262
cyclin B3	AB443443	802	270	0	270	34%	GE8SX9M02GFOJ4	271-368	704-802
cyclin E	AB378067	209	1276	617	659	611%	isotig01641	660-850	full

		209	1371	712	659	656%	isotig01640	660-850	full
		209	1446	787	659	692%	isotig01639	660-850	full
		209	3567	617	2896	1707%	isotig01638	2897-3087	full
		209	3608	712	2896	1726%	isotig01637	2897-3087	full
		209	3683	787	2896	1762%	isotig01636	2897-3087	full
DHHC-type zinc finger containing protein discs overgrown	AB378066	643	450	450	0	70%	isotig14108	353-1	29-381
	AB443442	287	3115	343	2772	1085%	isotig01394	344-630	full
		287	3211	439	2772	1119%	isotig01393	440-726	full
ecdysone receptor B1 EF1alpha	AB536932	828	435	0	435	53%	isotig14153	129-364	593-828
	AB583234	2029	101	101	0	5%	contig12129	101-604	1-155, 239-421, 490-660
		2029	116	0	116	6%	contig12130	117-505	1466-1543, 1627-1794, 1884-2029
elongation factor	DQ630923	717	432	0	432	60%	contig09678	433-779	371-717
		717	910	910	0	127%	contig09671	911-1048	1-138
Ena/VASP	AB378069	200	561	199	362	281%	isotig15279	363-515	25-177
enhancer of zeste	AB378079	431	1938	1423	515	450%	isotig05120	1424-1579, 1637-1854	1-218, 276-431
		431	2076	1423	653	482%	isotig05119	1424-1579, 1637-1854	1-218, 276-431
expanded	AB378099	648	182	0	182	28%	GFJY65E01B9CAF	1-284	346-629
		648	0	0	0	0%	GFJY65E01DAZCK	full	63-343
fasciclin-like protein fmr	DQ630929	768	1870	935	936	243%	isotig09432	936-1703	full
	AB461422	1854	42	42	0	2%	isotig06512	130-1053	88-1011
		1854	0	0	0	0%	isotig17262	full	1284-1825
epidermal growth factor receptor	AB300616	3807	0	0	0	0%	isotig12088	full	2369-3450
		3807	275	0	275	7%	isotig18881	276-456	3625-3807
GB1-cadherin	AB190295	4945	38	38	0	1%	isotig04828	1164-4545	1-3382

		4945	125	0	125	3%	isotig10276	126-1688	3383-4945
GB2-cadherin	AB190296	4096	203	203	0	5%	GFCP6CO01ER8W2	1-176	1-176
		4096	0	0	0	0%	FQTBZRY01B5B7K	full	875-1011
grainy head	AB378081	826	1244	0	1244	151%	isotig10851	1-225	579-803
		826	0	0	0	0%	FQTBZRY02GC2DO	full	7-226
Gug gene corepressor Atro 3'	AB378078	192	581	101	480	303%	isotig14567	102-166, 232-288	6-70, 136-192
Gug gene corepressor Atro 5'	AB378077	179	1671	1011	660	934%	isotig09993	1012-1151	17-156
hedgehog	AB044709	2963	0	0	0	0%	GE8SX9M01BZRKW	full	2142-2471
hexokinase	DQ630934	432	1539	769	770	356%	isotig09401	770-1201	full
hippo	AB378070	632	993	136	857	157%	isotig03128	137-768	full
		632	1131	136	995	179%	isotig03127	137-768	full
		632	321	136	185	51%	isotig03129	137-640	1-504
hunchback	AB120735	2672	0	0	0	0%	GFJY65E01C2FLA	full	2062-2295
		2672	0	0	0	0%	GE8SX9M02GCICC	full	2323-2669
inhibitor of apoptosis protein	AB378071	542	1253	168	1085	231%	isotig03633	1086-1628	full
		542	2241	1156	1085	413%	isotig03632	1086-1527	102-543
Insulin receptor	AB557977	386	865	431	434	224%	isotig04919	435-783	full
		386	3991	3558	434	1034%	isotig04918	435-783	full
kibra	DC445461	677	464	464	0	69%	isotig12669	567-971	103-507
		677	0	0	0	0%	isotig19618	full	10-560
		677	0	0	0	0%	isotig13198	full	1-669
		677	256	0	256	38%	isotig19193	332-285	601-648
merlin	AB378073	539	3525	3174	351	654%	isotig07940	3175-3712	full
mob as tumor suppressor	AB378072	381	1482	424	1058	389%	isotig09892	1059-1439	full
Musashi	AB459508	354	345	345	0	97%	GFJY65E02G1KQY	346-415	1-70
nitric oxide synthase	AB477987	3535	0	0	0	0%	GE8SX9M02FRE69	full	2233-2676
		3535	0	0	0	0%	GFCP6CO01EGMNE	full	2299-2671
		3535	0	0	0	0%	GFJY65E01DHCQJ	full	1751-2105

		3535	163	0	163	5%	GFJY65E02G0F7L	1-254	3272-3525
Notch	AB635585	2304	0	0	0	0%	isotig14599	full	374-1132
		2304	0	0	0	0%	isotig12243	full	1595-2300
		2304	0	0	0	0%	GE8SX9M01BNVPA	full	1145-1566
		2304	194	194	0	8%	GFCP6CO01B89FU	198-229	4-35
orthodenticle1	AB468156	720	599	0	599	83%	isotig12009	1-519	200-707
period	AB375516	3552	0	0	0	0%	isotig11839	436-1044	101-701
		3552	0	0	0	0%	GFCP6CO01E0OUQ	18-150, 509-538	1999-2155
phosphatase and tensin polycomb protein	AB370293	490	870	460	410	178%	isotig11178	411-775	full
	AB444104	1333	503	0	503	38%	isotig14622	504-761	1062-1319
		1333	0	0	0	0%	GE8SX9M02FSJ8W	full	789-1059
Ras association family member	AB443439	442	3221	133	3088	729%	isotig05452	3089-3530	full
		442	133	133	0	30%	isotig05453	305-172	1-134
s6k	AB557979	497	2738	491	2247	551%	isotig08277	2248-2744	full
S9 ribosomal protein	DQ630939	552	218	82	136	39%	isotig06773	82-626	full
		552	408	82	326	74%	isotig03301	82-398	1-420
salvador	AB378074	347	170	170	0	49%	GFJY65E01DC652	171-477	16-322
semaphorin 2a	EF036538	1306	203	0	203	16%	GFJY65E01CQDHA	1-49	1256-1303
sex combs reduced	AB194276	1015	0	0	0	0%	FQTBZRY02F97XW	4-34, 112-261	382-531, 609-639
Target of rapamycin	AB557078	269	230	12	218	86%	GFCO6CO02JFLMZ	13-277	full
tgf alpha (EGFR ligand)	HM106520	520	1476	0	1476	284%	isotig10026	1-323	194-520
timeless	AB548625.1	5795	0	0	0	0	isotig11684	full	3597-4790
		5795	0	0	0	0	isotig12618	full	4793-5784
		5795	0	0	0	0	isotig13095	full	2685-3596
		5795	0	0	0	0	isotig10108	1-722, 737-954, 1040-1758	622-1339, 1357-1574, 1660-2378
		5795	0	0	0	0	isotig15714	full	1-591
		5795	0	0	0	0	GFJY65E01EDCUE	full	2382-2663

<i>Ultrabithorax</i>	AB194278	790	0	0	0	0%	GFJY65E02JHD6P	361-58, 34-1	197-495,
		790	0	0	0	0%	GE8SX9M01BU9CB	1-135, 157-356	519-552, 362-496, 519-720
<i>vasa</i>	AB378065	1953	420	0	420	22%	isotig11874	1146-421	1228-1953
		1953	0	0	0	0%	isotig14543	full	439-1200
<i>Warts kinase</i>	AB300574	861	93	93	0	11%	isotig14894	641-1	19-659
<i>wingless</i>	AB044713	2298	0	0	0	0%	GE8SX9M01DIP5X	full	1610-2075
<i>yorkie</i>	AB378076	1021	403	403	0	39%	GFCP6CO03JM8ZY	404-507	1-102

Table S3

Transcription factors from the FlyTF database with putative orthologues identified in the *de novo* *G. bimaculatus* transcriptome.

AGO1	Camta	CG16778	CG4617	CG9305	Dcr-2	fd3F	jim
ab	cas	CG16779	CG4789	CG9416	Deaf1	Fen1	jing
abo	caz	CG16903	CG4882	CG9418	Dhc16F	fru	jumu
Acf1	cdc2	CG17209	CG5147	CG9425	Dip1	fry	kay
Ada2b	Cdk7	CG17829	CG5245	CG9437	dl	fs(1)h	Kdm4A
Alh	Cdk8	CG17912	CG5316	CG9705	DLP	ftz-f1	Kdm4B
alien	Cdk9	CG1832	CG5343	CG9817	dom	Gas41	ken
aop	cg	CG18619	CG5380	CG9890	Dp	GATAd	kin17
Arc42	CG10289	CG1965	CG5591	CG9932	Dp1	gce	king-tubby
ash1	CG10348	CG2712	CG5641	chinmo	DppIII	gcl	Kr-h1
ash2	CG10414	CG2790	CG5690	chm	Dref	gl	kto
Asx	CG10431	CG31211	CG5953	Chrac-14	Dsp1	gol	l(2)37Cg
Atac1	CG10565	CG31716	CG6129	Chrac-16	dys	grh	l(2)k10201
Atf6	CG10979	CG32121	CG6654	Chro	E(bx)	grn	l(2)NC136
aub	CG11414	CG32343	CG6686	cic	e(r)	gro	l(3)mbt
bab2	CG11456	CG3281	CG6701	Clk	e(y)1	grp	La
bap	CG11617	CG32830	CG6751	cnc	e(y)2	Gug	lack
Bap170	CG11710	CG32982	CG6765	Cog7	e(y)3	H	lds
Bap55	CG11876	CG3328	CG6769	CoRest	E(z)	h	lid
Bap60	CG12071	CG33695	CG6812	Cp190	E2f	hay	LIMK1
bbx	CG12162	CG33785	CG6905	crc	E2f2	Hcf	lin-52
bic	CG12236	CG33936	CG6907	CrebA	ear	hep	Lmpt
bigmax	CG12267	CG3407	CG7099	CrebB-17A	ecd	Hira	lola
Bin1	CG12299	CG34422	CG7339	CREG	EcR	HLH106	lolal
bip2	CG1233	CG3680	CG7368	crm	ECSIT	Hnf4	Mad
Bka	CG12370	CG3711	CG7556	croc	egg	HP1b	maf-S
bon	CG12769	CG3726	CG7785	crol	Eip74EF	hpo	mamo
br	CG13204	CG3735	CG7818	ct	Eip78C	Hr39	Mat1
brat	CG13458	CG3756	CG7839	CtBP	Eip93F	Hr4	Max
Brd8	CG13624	CG3815	CG7987	CTCF	EloA	Hr78	MBD-like
Brf	CG14200	CG3838	CG8152	CycC	Elongin-B	Hr96	MBD-R2
brk	CG14767	CG3909	CG8290	CycH	Elp3	Hsf	mbf1
brm	CG14962	CG40196	CG8359	CycT	emc	hth	Med
bs	CG15011	CG4042	CG8578	CYLD	ERR	lswi	MED1
BtbVII	CG15270	CG4404	CG8765	D12	Ets97D	ix	MED11
bun	CG15436	CG4553	CG8909	d4	ewg	Jarid2	MED14
Caf1	CG1620	CG4557	CG8924	dalao	exd	JIL-1	MED15

MED16	Mtp	phtf	Rpb8	Spt6	tld
MED17	mTTF	piwi	Rpd3	Ssdp	tna
MED18	mtTFB1	Pms2	Rpl1	Ssl1	Top2
MED20	mtTFB2	pnt	Rpl12	Ssrp	Top3alpha
MED21	mus201	polybromo	Rpl135	Stat92E	tou
MED22	mus308	Pop2	Rpl140	stc	tral
MED23	mus309	ppl	Rpl15	su(Hw)	Trax
MED24	Myb	pps	Rpl18	su(s)	Trf2
MED25	N	Psc	Rpl215	Su(var)2-10	Trl
MED27	Nap1	Psf2	Rpl33	Su(var)205	Trn-SR
MED28	NC2alpha	psq	Rpl1128	Su(var)3-9	trr
MED30	nej	pum	RpL40	Su(z)12	trsn
MED31	NELF-A	Pur-alpha	RpL7	sug	trx
MED4	NELF-B	put	sa	svp	ttk
MED6	Nelf-E	pygo	Sap30	Taf1	Tudor-SN
MED7	Nf-YA	pzg	Scamp	Taf10	Ubi-p63E
MED8	Nf-YB	r	Sce	Taf11	Unr
Meics	Nf-YC	Rab-RP4	Scm	Taf12	Usf
melt	Nipped-A	Rab1	Set2	Taf13	usp
MEP-1	Nipped-B	Rab10	Sfmbt	Taf2	Utx
Mes-4	nos	Rab11	sgg	Taf4	wash
Mes4	Not1	Rab2	sim	Taf5	wtb
MESR4	Nufip	Rab26	sima	Taf6	Xbp1
Met	opa	Rab27	simj	Taf8	XNP
Mi-2	Orc1	Rab35	Sin3A	tai	Xpd
mib1	Orc2	Rab8	Sir2	tara	yki
Mio	Orc5	Rbf	Sirt2	Tbp	YL-1
mip120	osa	Rel	Sirt4	tefu	yps
mip130	ovo	rept	Sirt6	TFAM	zfh1
mip40	p53	Rfx	Sirt7	Tfb1	zfh2
Mitf	pad	Rga	skd	Tfb4	Zpr1
Mlh1	Parg	rhea	Smox	TfIIA-L	
mod(mdg4)	Parp	rig	Smr	TfIIB	
mor	Pbp49	rl	sno	TfIIealpha	
MRG15	Pbp95	rn	Snr1	TfIIebeta	
mrn	Pc	rno	Sp1	TfIIFalpha	
mRpl28	Pcf11	row	spel1	TfIIFbeta	
mRpl55	Pcl	RpA-70	spen	TfIIS	
Msh6	peb	Rpb10	spn-A	TH1	
msl-3	pfk	Rpb11	spn-E	Thd1	
MTA1-like	ph-d	Rpb4	Spt3	Tif-IA	
MTF-1	ph-p	Rpb5	spt4	tim	
mtg	pho	Rpb7	Spt5	tkv	

Table S4

Selected signaling pathway genes identified in the *de novo* *G. bimaculatus* transcriptome.

Process	# Hits	Hit ID (I/S)	Length (range)	Query Gene	Transcriptome Sequence Name(s)
HEDGEHOG					
<i>CK1</i>	1	A	3248	<i>Ck1 alpha</i>	isotig08262
	2	A	3402-3498	<i>dco</i>	isotig01394, isotig01393
	1	A	2691	<i>gish</i>	isotig08729
<i>Cos2</i>	1	A	4125	<i>cos</i>	isotig07930
<i>Fused</i>	1	A	1624	<i>fu</i>	isotig10451
<i>TGFb</i>	1	A	1625	<i>gbb</i>	isotig07565
<i>GSK-3β</i>	2	A, S	367-483	<i>sgg</i>	GFJY65E02I1Z50, isotig18361
<i>Megalin</i>	1	A	2667	<i>Cg42611</i>	isotig08756
<i>Patched</i>	2	S	328-411	<i>ptc</i>	GFJY65E02I1VDN, GFJY65E01ALZ8M
<i>PKA</i>	1	A	4812	<i>Pka-C1</i>	isotig07789
<i>Smoothened</i>	2	A	705	<i>smo</i>	isotig13374, isotig15392
<i>Suppressor of fused</i>	1	A	2625	<i>Su(fu)</i>	isotig08905
<i>Slim b</i>	1	A	4768	<i>slmb</i>	isotig04954
JAK/STAT					
<i>AKT</i>	1	A	2629	<i>Akt1</i>	isotig08797
<i>Cb1</i>	1	A	486	<i>Cb1</i>	isotig18303
<i>CBP</i>	4	A, S	200-1501	<i>nej</i>	isotig17362, isotig05855, GE8SX9M02I88X1, isotig13864
<i>PIAS</i>	2	A	4065-4260	<i>Su(var)2-10</i>	isotig04583, isotig04582
<i>GRB</i>	1	A	2371	<i>drk</i>	isotig00085
<i>JAK</i>	1	A	2719	<i>hop</i>	isotig04276

<i>PI3K</i>	1	A	5218	<i>Pi3K21B</i>	isotig07744
	1	A	1976	<i>Pi3K92E</i>	isotig08270
<i>SHP2</i>	1	S	266	<i>csu</i>	GE8SX9M02G96K3
<i>SOCS</i>	1	A	2289	<i>Socs16D</i>	isotig09205
	1	A	3530	<i>Socs44A</i>	isotig08127
	2	A	2127-2190	<i>Socs36E</i>	isotig05318, isotig05317
<i>SOS</i>	5	A, S	170-1931	<i>Sos</i>	isotig09775, GFJY65E01CUEPZ, GFCP6CO02F9P6M, GAP9EXG04D7UG1, isotig14668
<i>Spred</i>	2	A	1189-3475	<i>Spred</i>	isotig05180, isotig05181
<i>STAM</i>	2	S	315-520	<i>Stam</i>	GFJY65E02GH370, GE8SX9M02IFMFR
<i>STAT</i>	1	A	2243	<i>Stat92E</i>	isotig03185
NOTCH					
<i>APH-1</i>	1	A	4738	<i>aph-1</i>	isotig04141
<i>CIR</i>	1	A	1088	<i>CG6843</i>	contig11433
<i>CtBP</i>	3	A, S	239-624	<i>CtBP</i>	isotig16142, GE8SX9M01EF4BJ, FQTBZRY01BYCPR
<i>Deltex</i>	2	A	1825-2309	<i>dx</i>	isotig09973, isotig09188
<i>disheveled</i>	2	A	2448-5763	<i>dsh</i>	isotig07449, isotig07448
<i>Groucho</i>	3	A, S	211-515	<i>gro</i>	isotig17698, GFJY65E01DLKWU, GFJY65E02GG7B9
<i>HATs</i>	4	A, S	200-1501	<i>nej</i>	isotig17362, isotig05855, GE8SX9M02I88X1, isotig13864
<i>HDAC</i>	1	A	2212	<i>Rpd3</i>	isotig09325
<i>Nicastrin</i>	4	A	766-2581	<i>nct</i>	isotig03085, isotig03084, isotig05814, isotig05814
<i>Notch</i>	4	A, S	423-2816	<i>Notch</i>	isotig14599, GE8SX9M01BNVPA, isotig12243, isotig08601
<i>Presenilin</i>	2	A	1999-3017	<i>Psn</i>	isotig03035, isotig03036
<i>PSE2</i>	1	A	864	<i>pen-2</i>	isotig13452
<i>SKIP</i>	3	A, S	338-2107	<i>Bx42</i>	isotig05493, isotig05494, GFJY65E02IALF1
<i>Tace</i>	1	A	3117	<i>Tace</i>	isotig08377
WNT					
<i>APC</i>	4	S	208-470	<i>Apc</i>	GFCP6CO02IKY6E, GFCP6CO01CGKAB, GFJY65E01EDMSG, GFJY65E01D5QKT
<i>Axin</i>	2	A	1769-2651	<i>Axn</i>	isotig00276, isotig08771

<i>beta-catenin</i>	1	A	3974	<i>arm</i>	isotig05341
<i>beta-TrCP</i>	1	A	4768	<i>slmb</i>	isotig04954
<i>CaMKII</i>	2	A	1262-2572	<i>CaMKII</i>	isotig05571, isotig05572
<i>CaN</i>	1	A	3292	<i>CanB2</i>	isotig05734
<i>CBP</i>	4	A, S	200-1501	<i>nej</i>	isotig17362, isotig05855, GE8SX9M02I88X1, isotig13864
<i>CK1</i>	1	A	3248	<i>Ck1 alpha</i>	isotig08262
	2	A	3402-3498	<i>dco</i>	isotig01394, isotig01393
	1	A	2691	<i>gish</i>	isotig08729
<i>CK2</i>	2	A	3799-4012	<i>Ckl beta</i>	isotig02546, isotig02545
<i>CtBP</i>	3	A, S	239-624	<i>CtBP</i>	isotig16142, GE8SX9M01EF4BJ, FQTBZRY01BYCPR
<i>Cul1</i>	4	A	3731-5498	<i>lin19</i>	isotig02458, isotig02457, isotig03607, isotig03606
<i>Daam1</i>	1	S	223	<i>DAAM</i>	FQTBZRY02HV44R
<i>disheveled</i>	2	A	2448-5763	<i>dsh</i>	isotig07449, isotig07448
<i>Ebi1</i>	1	A	2312	<i>ebi</i>	isotig09177
<i>GSK-3β</i>	2	A, S	367-483	<i>sgg</i>	GFJY65E02I1Z50, isotig18361
<i>Groucho</i>	3	A, S	211-515	<i>gro</i>	isotig17698, GFJY65E01DLKWU, GFJY65E02GG7B9
<i>JNK</i>	1	S	230	<i>bsk</i>	GFJY65E01CRM61
<i>LRP5/6</i>	2	S	259-493	<i>arr</i>	GFCP6CO01EVQLD, FQTBZRY02HHNYA
<i>NLK</i>	1	A	3303	<i>nmo</i>	isotig04244
<i>PKA</i>	1	A	4812	<i>Pka-C1</i>	isotig07789
<i>PKC</i>	1	A	4789	<i>Pkc53E</i>	isotig07795
<i>PLC</i>	1	S	329	<i>norpA</i>	GFJY65E01AO8M1
<i>PP2A</i>	3	A	1910-5172	<i>Pp2A-29B</i>	isotig02130, isotig02129, isotig09820
	1	A	1734	<i>mts</i>	isotig00164
<i>Proc</i>	1	A	1974	<i>por</i>	isotig09691
<i>Protein52</i>	1	A	1461	<i>pont</i>	contig15673
<i>PS-1</i>	2	A	1999-3017	<i>Psn</i>	isotig03035, isotig03036
<i>Rac</i>	1	A	2954	<i>Rac1</i>	isotig08497
<i>Rbx1</i>	2	S	459-480	<i>Roc1a</i>	GFCP6CO02GX3GB, GFCP6CO02I0JF4
<i>RhoA</i>	2	A	2482-3812	<i>Rho1</i>	isotig00258, isotig03933
<i>rhomboid-7</i>	1	A	3315	<i>rho-7</i>	isotig05079

<i>ROCK2</i>	8	A	853-4515	<i>rok</i>	isotig01612, isotig01613, isotig01614, isotig01615, isotig01616, isotig01617, isotig06106, isotig06107
<i>Siah-1</i>	2	A	1698-2386	<i>sina</i>	isotig09073, isotig10251
	8	A	3812-4068	<i>sinah</i>	isotig00589, isotig00588, isotig00587, isotig00586, isotig00585, isotig00584, isotig00583, isotig00582
<i>SIP</i>	1	A	587	<i>CG3226</i>	contig15490
<i>Skp1</i>	1	A	951	<i>skpF</i>	isotig12819
<i>SMAD3</i>	1	S	332	<i>Smox</i>	GFCP6CO01DS40Z
<i>SMAD4</i>	1	A	729	<i>Med</i>	isotig15042
<i>Stbm</i>	1	A	2916	<i>Vang</i>	isotig08532
<i>Wif-1</i>	2	A	1568-1581	<i>shf</i>	isotig02624, isotig02623
TGF-BETA					
<i>ActivinRl</i>	1	A	2267	<i>babo</i>	isotig09236
<i>Cul1</i>	4	A	3731-5498	<i>lin19</i>	isotig02458, isotig02457, isotig03607, isotig03606
<i>DP1</i>	1	A	3452	<i>tfdp1a</i>	isotig08163
<i>E2F4/5</i>	4	A	1708-1929	<i>e2f4</i>	isotig00805, isotig00806, isotig00807, isotig00808
<i>ERK</i>	1	A	799	<i>rl</i>	isotig14164
<i>Id</i>	1	S	201	<i>emc</i>	FQTBZRY02G5SHM
<i>TGFb</i>	1	A	1625	<i>gbb</i>	isotig07565
<i>p107</i>	2	A	6434-6542	<i>Rbf</i>	isotig04489, isotig04488
<i>p300</i>	4	A, S	200-1501	<i>nej</i>	isotig17362, isotig05855, GE8SX9M02I88X1, isotig13864
<i>p70S6K</i>	1	A	3234	<i>S6K</i>	isotig08277
<i>PP2A</i>	3	A	1910-5172	<i>Pp2A-29B</i>	isotig02130, isotig02129, isotig09820
	1	A	1734	<i>mts</i>	isotig00164
<i>Rbx1</i>	2	S	459-480	<i>Roc1a</i>	GFCP6CO02GX3GB, GFCP6CO02I0JF4
<i>RhoA</i>	2	A	2482-3812	<i>Rho1</i>	isotig00258, isotig03933
<i>ROCK1</i>	8	A	853-4515	<i>rok</i>	isotig01612, isotig01613, isotig01614, isotig01615, isotig01616, isotig01617, isotig06106, isotig06107
<i>SARA</i>	1	A	2592	<i>Sara</i>	isotig08835
<i>Skp1</i>	1	A	951	<i>skpF</i>	isotig12819
<i>Smad1/5/8</i>	1	A	2120	<i>Mad</i>	isotig09444
<i>Smad2/3</i>	1	S	332	<i>Smox</i>	GFCP6CO01DS40Z

<i>Smad4</i>	1	A	729	<i>Med</i>	isotig15042
<i>Smurf1/2</i>	1	A	4308	<i>lack</i>	isotig07879
MAPK					
<i>Boss</i>	1	A	3134	<i>boss</i>	isotig08354
<i>Csw</i>	1	S	266	<i>csw</i>	GE8SX9M02G96K3
<i>Drk</i>	1	A	2371	<i>drk</i>	isotig00085
<i>Dsor1</i>	1	A	3545	<i>Dsor1</i>	isotig08121
<i>Egfr</i>	1	A	1099	<i>Egfr</i>	isotig12088
<i>Gap1</i>	2	S	280-358	<i>Gap1</i>	GE8SX9M02HTUD8, GFJY65E01B2IBY
<i>Phl</i>	1	A	4282	<i>phl</i>	isotig07892
<i>Pointed</i>	1	S	314	<i>pnt</i>	GFCP6CO01CJJKD
<i>Ras85D</i>	2	A	2078-2467	<i>Ras85D</i>	isotig09494, isotig08979
<i>Rolled</i>	1	A	799	<i>rl</i>	isotig14164
<i>Sos</i>	5	A, S	170-1931	<i>Sos</i>	isotig09775, GFJY65E01CUEPZ, GFCP6CO02F9P6M, GAP9EXG04D7UG1, isotig14668
<i>Ts1</i>	1	S	174	<i>ts1</i>	GFCP6CO02G92YK
<i>Yan</i>	1	A	4007	<i>aop</i>	isotig07960
HIPPO					
<i>cyclinE</i>	6	A	1521-3799	<i>CycE</i>	isotig01638, isotig01637, isotig01636, isotig01641, isotig01640, isotig01639
<i>Dco</i>	2	A	3402-3498	<i>dco</i>	isotig01394, isotig01393
<i>diap1</i>	2	A	1796-2785	<i>th</i>	isotig03633, isotig03632
<i>Expanded</i>	2	S	306-466	<i>ex</i>	GFJY65E01B9CAF, GFJY65E01DAZCK
<i>Fat</i>	1	A	716	<i>ft</i>	isotig15250
<i>Hippo</i>	3	A	953-1763	<i>hpo</i>	isotig03127, isotig03128, isotig03129
<i>homothorax</i>	5	S	153-234	<i>hth</i>	FQTBZRY01D5WGD, FQTBZRY01AYECO, FQTBZRY02GY7HW, FQTBZRY02JLK81, FQTBZRY02FHQV8
<i>Kibra</i>	4	A	365-974	<i>kibra</i>	isotig12669, isotig19618, isotig19193, isotig13198
<i>Merlin</i>	1	A	1313	<i>Mer</i>	isotig11307
<i>Mob as tumor suppressor</i>	1	A	1862	<i>mats</i>	isotig09892
<i>Salvador</i>	1	S	479	<i>sav</i>	GFJY65E01DC652

<i>Warts</i>	1	S	300	<i>wtS</i>	GFJY65E01AT7SH
<i>yorkie</i>	1	S	507	<i>yki</i>	GFCP6CO02JM8ZY

Table S5

Selected gametogenesis genes identified in the *de novo* *G. bimaculatus* transcriptome

Process	# Hits	Hit ID (A/S)	Length (range)	Query Gene	Transcriptome Sequence Names
SPERMATOGENESIS¹					
TRANSCRIPTION FACTORS					
<i>Enhancer of bithorax</i>	1	A	4242	<i>E(bx)</i>	contig15318
<i>eyes absent</i>	1	S	401	<i>eya</i>	GFJY65E01EO7KL
<i>Heat shock factor</i>	4	A	3119-3268	<i>Hsf</i>	isotig01705, isotig01704, isotig01703, isotig01702
<i>maleless</i>	1	A	3818	<i>mle</i>	isotig05146
<i>MBD-like</i>	7	A	694-1211	<i>MBD-like</i>	isotig01061, isotig01060, isotig01064, isotig01063, isotig01062, isotig01066, isotig01065
<i>Myb oncogene-like</i>	1	A	3771	<i>Myb</i>	isotig08042
<i>Rfx</i>	1	A	1001	<i>Rfx</i>	isotig12547
<i>TATA box binding protein-related factor 2</i>	3	A, S	399-3469	<i>Trf2</i>	GFCP6CO01B8937, isotig01886, isotig01885
CYTOSKELETON					
<i>Adenomatous popylosis coli tumor suppressor homolog (APC)</i>	4	S	208-470	<i>Apc</i>	GFCP6CO02IKY6E, GFCP6CO01CGKAB, GFJY65E01EDMSG, GFJY65E01D5QKT
<i>Adenomatous popylosis coli tumor suppressor homolog 2</i>	2	S	315-470	<i>Apc2</i>	GFCP6CO02IKY6E, GFJY65E01D5QKT
<i>beta tubulin</i>	2	A	546-950	<i>Btub56D</i>	contig00262, contig00455

¹ Although we did not include cDNA derived from adult testes in our sequencing libraries, we nonetheless chose to perform manual annotation of genes known to be involved in *D. melanogaster* spermatogenesis since the creation of the testis germ line stem cell niche takes place during embryogenesis in *D. melanogaster* (Aboïm AN (1945) Développement embryonnaire et post-embryonnaire des gonades normales et agamétiques de *Drosophila melanogaster*. Revue Suisse de Zoologie 3: 53-154; Le Bras S, Van Doren M (2006) Development of the male germline stem cell niche in *Drosophila*. Developmental Biology 294: 92-103.) and in orthopterans 3. Nelsen OE (1931) Life cycle, sex differentiation, and testis development in *Melanoplus differentialis* (Acrididae, Orthoptera). Journal of Morphology 51: 467-525.)

<i>cortactin</i>	1	A	1147	<i>Cortactin</i>	isotig11852
<i>diaphanous</i>	1	S	237	<i>dia</i>	FQTBZRY01CIL7E
<i>jaguar</i>	3	A	958-2609	<i>jar</i>	isotig12791, isotig12012, isotig08822
<i>Kinesin like protein at 61F</i>	3	A	2102-3639	<i>Klp61F</i>	isotig01563, isotig01564, isotig01565
<i>Myosin 31DF</i>	2	A	1018-1306	<i>Myo31DF</i>	isotig11312, isotig12459
<i>peanut</i>	1	A	1957	<i>pnut</i>	isotig09723
<i>Rac1</i>	1	A	2954	<i>Rac1</i>	isotig08497
<i>Spectrin 1</i>	4	A	409-2155	<i>α-Spec</i>	isotig09397, isotig10052, isotig15468, isotig19330
<i>spindle assembly abnormal 6</i>	1	A	2655	<i>sas-6</i>	isotig05533
<i>subito</i>	1	A	2615	<i>sub</i>	contig14686
<i>twinstar</i>	2	A	513-2077	<i>tsr</i>	isotig00493, isotig00494
<i>zipper</i>	2	A	3077-3958	<i>zip</i>	isotig05158, isotig08407
OTHER PROCESSES IN SPERMATOGENESIS					
<i>armitage</i>	1	A	4095	<i>armi</i>	isotig07934
<i>asterless</i>	1	A	3788	<i>asl</i>	isotig08040
<i>aubergine</i>	2	A	2674-2784	<i>aub</i>	isotig07461, isotig07462
<i>boule</i>	1	S	203	<i>bol</i>	GFJY65E01B4FFK
<i>bride of sevenless</i>	1	A	3134	<i>boss</i>	isotig08354
<i>Btk family kinase at 29A</i>	2	A	915-1545	<i>Btk29A</i>	isotig06869, isotig10647
<i>Bub1-related kinase</i>	1	A	4209	<i>BubR1</i>	isotig07912
<i>Calmodulin</i>	3	A	1591-1698	<i>Cam</i>	isotig00266, isotig00265, isotig00264
<i>capsuleen</i>	2	A	3725-3816	<i>csul</i>	isotig01229, isotig01228
<i>cdc2</i>	1	A	2078	<i>cdc2</i>	isotig03292
<i>courtless</i>	1	A	1123	<i>crl</i>	isotig11993
<i>Cyclin A</i>	1	A	3049	<i>CycA</i>	isotig03226
<i>Cytochrome c proximal</i>	1	A	636	<i>Cyt-c-p</i>	contig10573
<i>Cytochrome c distal</i>	1	A	778	<i>Cyt-c-d</i>	isotig14404
<i>Dynamin related protein 1</i>	4	A	908-3502	<i>Drp1</i>	isotig13131, isotig01328, isotig01327, isotig01326
<i>effete</i>	2	A	3342-4080	<i>eff</i>	isotig01782, isotig01780
<i>Fmr1</i>	2	A	1038-1053	<i>Fmr1</i>	isotig06512, isotig06513
<i>Fps oncogene analog</i>	1	A	328	<i>Fps85D</i>	isotig19747
<i>glass bottom boat</i>	1	A	1625	<i>gbb</i>	isotig07565

<i>gilgamesh</i>	1	A	2691	<i>gish</i>	isotig08729
<i>hephaestus</i>	2	S	201-359	<i>heph</i>	GE8SX9M01A0TGF, FQTBZRY02F9D2F
<i>Ice</i>	2	A	1620-1800	<i>Ice</i>	isotig04366, isotig10455
<i>karyopherin α1</i>	1	A	1309	<i>Kap-α1</i>	isotig11303
<i>loquacious</i>	1	A	2867	<i>loqs</i>	isotig02873
<i>Microcephalin</i>	2	A	3384-4822	<i>MCPH1</i>	isotig04588, isotig04589
<i>Myt1</i>	1	A	3433	<i>Myt1</i>	isotig04225
<i>Nedd2-like caspase</i>	1	A	2900	<i>Nc</i>	isotig03487
<i>parkin</i>	2	A	3339-3502	<i>park</i>	isotig04723, isotig04722
<i>pavarotti</i>	2	A	2221-2661	<i>pav</i>	isotig03050, isotig03049
<i>pelota</i>	1	A	922	<i>pelo</i>	contig17247
<i>piwi</i>	1	A	1277	<i>piwi</i>	isotig11428
<i>pole hole</i>	1	A	4282	<i>phl</i>	isotig07892
<i>punt</i>	1	S	441	<i>put</i>	GE8SX9M01B9MGK
<i>Rab-protein 11</i>	1	A	2448	<i>Rab11</i>	isotig00835
<i>Rheb</i>	1	A	953	<i>Rheb</i>	contig21414
<i>shotgun</i>	1	A	4583	<i>shg</i>	isotig04828
<i>shut down</i>	2	A	2449-3029	<i>shu</i>	isotig04931, isotig04930
<i>string</i>	1	A	911	<i>stg</i>	isotig13103
<i>Syntaxin 5</i>	3	A	2683-3493	<i>Syx5</i>	isotig01824, isotig01823, isotig01825
<i>transformer 2</i>	1	A	836	<i>tra2</i>	contig12123
<i>terribly reduced optic lobes</i>	1	A	690	<i>trol</i>	isotig15574
<i>uncoordinated</i>	1	A	2116	<i>unc</i>	isotig09457
<i>vav</i>	1	A	2068	<i>vav</i>	isotig09529
<i>ypsilon schachtel</i>	1	A	2601	<i>yps</i>	isotig03079
OOGENESIS					
MAINTENANCE AND DIVISION OF GERM LINE STEM CELLS					
<i>armadillo</i>	1	A	3974	<i>arm</i>	isotig05341
<i>Axin</i>	2	A	1769-2651	<i>Axn</i>	isotig00276, isotig08771
<i>Dicer-1</i>	1	A	2177	<i>Dcr-1</i>	isotig09376
<i>dishevelled</i>	2	A	2448-5763	<i>dsh</i>	isotig07449, isotig07448
<i>effete</i>	2	A	3342-4080	<i>eff</i>	isotig01782, isotig01780

<i>fused</i>	1	A	1624	<i>fu</i>	isotig10451
					GFJY65E01C8HCB, GE8SX9M01D9LON, GFJY65E01EPKW2, GE8SX9M01D913W, FQTBZRY01EKMIL, GE8SX9M01AEJPJ, GFCP6CO01BN88A, GE8SX9M01ASUJ7, GFJY65E02HJ33N, isotig18880, GFCP6CO02F8AKG, GFCP6CO02GAOJB, isotig07261, GFCP6CO01AQ9N2, FQTBZRY02J3ED4, FQTBZRY01DAFOD
<i>karst</i>	16	A, S	140-568	<i>kst</i>	
<i>loquacious</i>	1	A	2867	<i>loqs</i>	isotig02873
<i>ovarian tumor</i>	2	A	2393-2483	<i>out</i>	isotig05114, isotig05113
<i>pelota</i>	1	A	922	<i>pelo</i>	contig17247
<i>piwi</i>	1	A	1277	<i>piwi</i>	isotig11428
<i>pumilio</i>	3	A, S	412-624	<i>pum</i>	isotig04477, isotig04476, GFJY65E02G1R75
<i>sans fille</i>	1	A	1511	<i>snf</i>	isotig10698
<i>shaggy</i>	1	A	483	<i>sgg</i>	isotig18361
<i>shavenbaby</i>	1	A	795	<i>ovo</i>	isotig14222
<i>shut down</i>	2	A	2449-3029	<i>shu</i>	isotig04931, isotig04930
<i>vasa</i>	2	A	765-1146	<i>vas</i>	isotig14543, isotig11874
OOCYTE DETERMINATION AND FORMATION OF AP AXIS					
<i>4EHP</i>	1	A	1414	<i>4EHP</i>	isotig01556
<i>alpha Spectrin</i>	4	A	409-2155	<i>α-Spec</i>	isotig09397, isotig10052, isotig15468, isotig19330
<i>beta-Tubulin at 56D</i>	2	A	546-950	<i>Btub56D</i>	contig00262, contig00455
<i>Bicaudal C</i>	2	A	854-1435	<i>BicC</i>	isotig06390, isotig06389
<i>Bicaudal D</i>	2	A	687-1014	<i>BicD</i>	isotig12488, isotig15621
<i>cAMP-dependent protein kinase 1</i>	1	A	4812	<i>Pka-C1</i>	isotig07789
<i>COP9 complex homolog subunit 5</i>	2	A	1032-1284	<i>CSN5</i>	contig13654, isotig11391
<i>cornichon</i>	1	A	1733	<i>cni</i>	isotig05694
					isotig15021, isotig12385, isotig18811, GFJY65E02JTGDA, GFJY65E01CXFIZ, isotig13703, isotig10229, isotig10644
<i>Dynein heavy chain 64C</i>	8	A, S	344-1706	<i>Dhc64C</i>	
<i>Dystroglycan</i>	2	S	293-342	<i>Dg</i>	GFCP6CO01C30LP, GFCP6CO01BPUA2
<i>egalitarian</i>	2	A	878-1634	<i>egl</i>	isotig13386, isotig10415
<i>egghead</i>	1	A	796	<i>egh</i>	isotig14205
<i>exuperantia</i>	2	A	3152-3225	<i>exu</i>	isotig04764, isotig04765

<i>Helicase at 25E</i>	2	S	277-341	<i>Hel25E</i>	GFJY65E01EGNY3, GE8SX9M01BJ16P
<i>hu-li tai shao</i>	6	A	2255-2885	<i>hts</i>	isotig01647, isotig01646, isotig01645, isotig01644, isotig01643, isotig01642
<i>Kinesin heavy chain</i>	2	A	3918-7009	<i>Khc</i>	isotig04492, isotig04493
<i>licorne</i>	1	A	2845	<i>lic</i>	contig18303
<i>lkb1</i>	2	A	3048-3216	<i>lkb1</i>	isotig01200, isotig01199
<i>maelstrom</i>	1	A	2668	<i>mael</i>	isotig06013
<i>okra</i>	1	A	1794	<i>okr</i>	isotig10034
<i>par-1</i>	1	A	889	<i>par-1</i>	isotig07610
<i>par-6</i>	1	A	3994	<i>par-6</i>	isotig07961
<i>pipsqueak</i>	1	A	1991	<i>Rab-6</i>	isotig09661
<i>tudor</i>	1	A	3025	<i>spn-E</i>	contig00220
FORMATION OF DV AXIS					
<i>cappuccino</i>	2	A	817-866	<i>capu</i>	isotig06798, isotig06799
<i>orb</i>	1	A	4765	<i>orb</i>	isotig00462
<i>pipe</i>	1	A	6608	<i>pip</i>	isotig07697
<i>squid</i>	1	A	1546	<i>sqd</i>	isotig00544
<i>trailer hitch</i>	2	A	263-493	<i>tral</i>	isotig18126, isotig07398
ACTING EARLY IN FOLLICLE CELLS (DORSAL GROUP)					
<i>big brain</i>	3	S	298-515	<i>bib</i>	GE8SX9M01BXNN0, GFJY65E01CFBEX, GFCP6CO01EV5QZ
<i>bunched</i>	1	A	869	<i>bun</i>	isotig13467
<i>Chorion factor 2</i>	1	S	147	<i>Cf2</i>	GFCP6CO01DSOAR
<i>corkscrew</i>	1	S	266	<i>csw</i>	GE8SX9M02G96K3
<i>dodo</i>	2	A	1975-1994	<i>dod</i>	isotig05499, isotig05500
<i>broad</i>	1	A	904	<i>br</i>	isotig13160
<i>torpedo</i>	1	A	1099	<i>Egfr</i>	isotig12088
<i>Ets at 97D</i>	1	A	2149	<i>Ets97D</i>	isotig05797
<i>kibra ortholog</i>	1	A	974	<i>kibra</i>	isotig12669
<i>mago nashi</i>	1	A	1021	<i>mago</i>	isotig12375
<i>Notch</i>	4	A, S	423-2816	<i>Notch</i>	isotig14599, GE8SX9M01BNVPA, isotig12243, isotig08601
<i>pointed</i>	1	S	314	<i>pnt</i>	GFCP6CO01CJJKD
<i>Rac1</i>	1	A	2954	<i>Rac1</i>	isotig08497

<i>Ras oncogene at 85D</i>	2	A	2078-2467	<i>Ras85D</i>	isotig09494, isotig08979
<i>rolled</i>	1	A	799	<i>rl</i>	isotig14164
<i>singed</i>	1	S	239	<i>sn</i>	GE8SX9M01EZ3K3
TERMINAL GENES					
<i>SHC-adaptor protein</i>	2	A	2374-2640	<i>Shc</i>	isotig05081, isotig05082
<i>torso-like</i>	1	S	174	<i>ts1</i>	GFCP6CO02G92YK
LIGANDS, RECEPTORS & EFFECTORS					
<i>hopscotch</i>	1	A	2719	<i>hop</i>	isotig04276
<i>Keren</i>	1	A	1803	<i>Krn</i>	isotig10026
<i>kugelei</i>	1	A	729	<i>kug</i>	isotig15037
<i>Medea</i>	1	A	729	<i>Med</i>	isotig15042
<i>Mothers against dpp</i>	1	A	2120	<i>Mad</i>	isotig09444
<i>Protein tyrosine phosphatase 69D</i>	2	A, S	471-1475	<i>Ptp69D</i>	isotig10837, GFCP6CO02HK6UL
<i>punt</i>	1	S	441	<i>put</i>	GE8SX9M01B9MGK
<i>saxophone</i>	1	A	4561	<i>sax</i>	isotig07822
<i>shotgun</i>	1	A	4583	<i>shg</i>	isotig04828
<i>Star</i>	1	A	4011	<i>S</i>	isotig07955
<i>STAT</i>	1	A	2243	<i>Stat92E</i>	isotig03185
GENES AFFECTING CYTOSKELETON					
<i>adnormal spindle</i>	1	A	6563	<i>asp</i>	isotig07699
<i>alpha actinin</i>	1	A	2837	<i>Actn</i>	isotig08592
<i>Btk family kinase at 29A</i>	2	A	915-1545	<i>Btk29A</i>	isotig06869, isotig10647
<i>capulet</i>	1	A	3379	<i>capt</i>	isotig04236
<i>Cdc42</i>	1	A	2958	<i>Cdc42</i>	isotig03915
<i>Ced-12</i>	1	A	3012	<i>Ced-12</i>	isotig08450
<i>chromosome bows</i>	1	A	1067	<i>chb</i>	isotig12228
<i>sticky</i>	1	A	3121	<i>sti</i>	isotig08364
<i>Cortactin</i>	1	A	1147	<i>Cortactin</i>	isotig11852
<i>diaphanous</i>	2	S	237-429	<i>dia</i>	FQTBZRY01CIL7E, GFJY65E01CBNCA
<i>qenghis khan</i>	1	A	2408	<i>gek</i>	isotig09046
<i>Jaquar</i>	3	A	958-2609	<i>jar</i>	isotig12791, isotig12012 , isotig08822
<i>kette</i>	1	A	5316	<i>Hem</i>	isotig07736
<i>Kinesin associated protein 3</i>	1	A	3027	<i>Kap3</i>	contig12721
<i>klarsicht</i>	1	A	1805	<i>klar</i>	isotig10023

<i>Lamin</i>	1	A	1757	<i>Lam</i>	contig17155
<i>Lissencephaly</i>	1	A	4309	<i>Lis-1</i>	isotig02186
<i>mushrom body defect</i>	1	A	2026	<i>mud</i>	contig12641
					GFJY65E01DWNJO,GFJY65E01DDZGX, isotig00295, isotig00293, GE8SX9M02JKH1C , GE8SX9M01EGTZF , GFCP6C001CRFSJ, GE8SX9M02F0G9A
<i>rho-type guanine exchange factor</i>	8	A, S	234-1847	<i>rtGEF</i>	
<i>short stop</i>	3	A	673-1571	<i>shot</i>	isotig13049, isotig10577, isotig15743
<i>spaghetti squash</i>	3	A	616-1053	<i>sqh</i>	contig15080, isotig00107, isotig00106
<i>Src oncogene at 42A</i>	1	A	1787	<i>Src42A</i>	isotig04219
<i>subito</i>	1	A	2615	<i>sub</i>	contig14686
<i>Suppressor of profilin 2</i>	1	A	1695	<i>Sop2</i>	isotig06657
<i>twinstar</i>	2	A	513-2077	<i>tsr</i>	isotig00493, isotig00494
<i>washout</i>	1	A	608	<i>wash</i>	isotig07224
<i>zipper</i>	1	A	3958	<i>zip</i>	isotig05158
OTHER GENES INVOLVED IN OOGENESIS					
<i>altered disjunction</i>	3	A	3259-3423	<i>ald</i>	isotig03616, isotig03615, isotig03614
<i>archipelago</i>	1	A	4393	<i>ago</i>	isotig00333
<i>chiffon</i>	1	A	3144	<i>chif</i>	isotig08349
<i>Cyclin-dependent kinase 7</i>	1	A	2272	<i>Cdk7</i>	isotig02269
<i>Cyclin-dependent kinase subunit30A</i>	1	A	1011	<i>Cks30A</i>	isotig06131
<i>Cyclin E</i>	6	A	1521-3799	<i>CycE</i>	isotig01638, isotig01637, isotig01636, isotig01641, isotig01640, isotig01639
<i>double parked</i>	1	A	5242	<i>dup</i>	isotig07741
<i>E2F transcription factor</i>	1	A	918	<i>E2f</i>	isotig13069
<i>geminin</i>	1	A	921	<i>geminin</i>	isotig04440
<i>imaginal discs arrested</i>	1	A	5406	<i>ida</i>	isotig07735
<i>loki</i>	1	A	2488	<i>lok</i>	isotig05744
<i>meiotic 41</i>	1	A	1228	<i>mei-41</i>	isotig11599
<i>Microcephalin</i>	2	A	3384-4822	<i>MCPH1</i>	isotig04588, isotig04589
<i>morula</i>	2	A	1648-1877	<i>mr</i>	isotig09875, isotig10364
<i>mutagen-sensitive 209</i>	1	A	1396	<i>mus209</i>	isotig00238
<i>Myb oncogene-like</i>	1	A	3771	<i>Myb</i>	isotig08042
<i>Myt1</i>	1	A	3433	<i>Myt1</i>	isotig04225
<i>pitchoune</i>	1	A	3126	<i>pit</i>	isotig04252
<i>sarah</i>	1	A	3505	<i>sra</i>	isotig08135
<i>twins</i>	2	A	2511-3747	<i>tw</i>	isotig04782, isotig04783
<i>abstrakt</i>	1	A	1429	<i>abs</i>	isotig10965

<i>anterior open</i>	1	A	4007	<i>aop</i>	isotig07960
<i>aubergine</i>	2	A	2674-2784	<i>aub</i>	isotig07461, isotig07462
<i>Autophagy-specific gene 1</i>	1	A	1467	<i>Atg1</i>	contig16688
<i>basket</i>	1	S	230	<i>bsk</i>	GFJY65E01CRM61
<i>blistered</i>	2	S	233-241	<i>bs</i>	GE8SX9M02GMAG7, GFJY65E02FKS1J
<i>brainiac</i>	1	A	1870	<i>brn</i>	isotig09883
<i>Bruce</i>	1	A	4923	<i>Bruce</i>	isotig07779
<i>capsuleen</i>	2	A	3725-3816	<i>csul</i>	isotig01229, isotig01228
<i>Calmodulin-binding protein related to a Rab3 GDP/GTP exchange protein</i>	2	A	1233-2207	<i>Crag</i>	isotig07677, isotig09331
<i>combgap</i>	1	A	3604	<i>cg</i>	isotig08105
<i>Cyclic-AMP response element binding protein A</i>	1	A	3290	<i>CrebA</i>	isotig08237
<i>C-terminal binding protein</i>	3	A, S	239-624	<i>CtBP</i>	isotig16142, GE8SX9M01EF4BJ, FQTBZRY01BYCPR
<i>cut</i>	1	S	247	<i>ct</i>	FQTBZRY02GRN27
<i>Death related ced-3/Nedd2-like protein</i>	1	A	2732	<i>Dredd</i>	isotig08688
<i>Ecdysone-induced protein 63E</i>	2	A	3994-4024	<i>Eip63E</i>	isotig02121, isotig02120
<i>ecdysoneless</i>	4	A, S	372-3069	<i>ecd</i>	isotig17485, isotig19531, GFCEP6CO01D3B2B, isotig08412
<i>eggless</i>	2	A	2986-3019	<i>egg</i>	isotig04831, isotig04830
<i>extra macrochaetae</i>	1	S	201	<i>emc</i>	FQTBZRY02G5SHM
<i>fat facets</i>	5	A	1816-3259	<i>faf</i>	isotig01188, isotig01187, isotig01186, isotig01185, isotig01184
<i>fruitless</i>	2	A	1313-1618	<i>fru</i>	isotig06010, isotig06009
<i>G protein-coupled receptor kinase 2</i>	1	A	1632	<i>Gprk2</i>	isotig00416
<i>G protein α 47A</i>	1	A	2901	<i>G-α47A</i>	isotig05513
<i>poly U binding factor 68kD</i>	2	A	3724-3736	<i>pUf68</i>	isotig01566, isotig01567
<i>Heat shock factor</i>	4	A	3119-3268	<i>Hsf</i>	isotig01705, isotig01704, isotig01703, isotig01702
<i>Heat-shock-protein-70</i>	3	A	2209-2595	<i>Hsp70</i>	isotig09115, isotig00207, isotig00208
<i>Hepatocyte growth factor regulated tyrosine kinase substrate</i>	4	A, S	344-583	<i>Hrs</i>	isotig16755, GE8SX9M02I536Z, GE8SX9M01EC494, GE8SX9M01BQDGK
<i>hephaestus</i>	2	S	201-359	<i>heph</i>	GE8SX9M01A0TGF, FQTBZRY02F9D2F
<i>Ice</i>	2	A	1620-1800	<i>Ice</i>	isotig04366, isotig10455
<i>jing</i>	2	S	366-427	<i>jing</i>	GFCEP6CO01CJPNC, GE8SX9M02FPBO4
<i>jumeau</i>	1	A	3251	<i>jumu</i>	isotig08268
<i>leonardo</i>	2	A	3053-3220	<i>14-3-3ζ</i>	isotig03712, isotig03711
<i>lethal (2) giant larvae</i>	2	A	1879-2573	<i>l(2)gl</i>	contig15364, contig15365
<i>Lipid storage droplet-2</i>	1	A	1861	<i>Lsd-2</i>	isotig06100
<i>Liprin-a</i>	3	A	982-1158	<i>Liprin-a</i>	isotig03903, isotig03902, isotig03901

<i>maternal expression at 31B</i>	1	A	3408	<i>me31B</i>	isotig00511
<i>Merlin</i>	1	A	4062	<i>Mer</i>	isotig07940
<i>Methoprene-tolerant</i>	1	S	420	<i>Met</i>	GFCP6CO01DP0NH
<i>microtubule star</i>	1	A	1734	<i>mts</i>	isotig00164
<i>mini spindles</i>	2	A	2315-4784	<i>msps</i>	isotig07797, isotig09181
<i>misshapen</i>	1	S	327	<i>msn</i>	GE8SX9M01DYBJ8
<i>moira</i>	1	A	2758	<i>mor</i>	isotig08664
<i>Nedd2-like caspase</i>	1	A	2900	<i>Nc</i>	isotig03487
<i>nicastatin</i>	2	A	2050-1886	<i>nct</i>	isotig03085, isotig03084
<i>Niemann-Pick type C-2a</i>	1	A	1095	<i>Npc2a</i>	contig15402
<i>Nucleolar protein at 60B</i>	1	A	878	<i>Nop60 B</i>	contig09572
<i>O-fucosyltransferase 1</i>	1	A	3328	<i>O-fut1</i>	isotig08223
<i>Ornithine decarboxylase antizyme</i>	1	A	2621	<i>Oda</i>	isotig08802
<i>PDGF- and VEGF-receptor related</i>	1	A	2602	<i>Pvr</i>	isotig08816
<i>pollux</i>	1	A	3676	<i>plx</i>	isotig08081
<i>polyhomeotic distal</i>	1	A	1608	<i>ph-d</i>	isotig10480
<i>polyhomeotic proximal</i>	1	S	514	<i>ph-p</i>	GFCP6CO02H0634
<i>Presenilin</i>	2	A	1999-3017	<i>Psn</i>	isotig03035, isotig03036
<i>Rab-protein 5</i>	2	A	3300-3532	<i>Rab5</i>	isotig02948, isotig02947
<i>Rab-protein 11</i>	1	A	2448	<i>Rab11</i>	isotig00835
<i>rotund</i>	1	S	186	<i>rn</i>	FQTBZRY02G28N8
<i>scribbled</i>	2	A, S	427-696	<i>scrib</i>	isotig15514, GE8SX9M01C2HSN
<i>skittles</i>	1	S	311	<i>sktl</i>	GE8SX9M01DH70U
<i>SNF1A/AMP-activated protein kinase</i>	1	A	2566	<i>SNF1A</i>	isotig05865
<i>Snf5-related 1</i>	1	A	3127	<i>Snr1</i>	isotig08358
<i>spinster</i>	1	A	3143	<i>spin</i>	isotig00443
<i>SH2 ankyrin repeat kinase</i>	1	A	3892	<i>shark</i>	isotig07997
<i>strawberry notch</i>	2	A	2330-2459	<i>sno</i>	isotig09159, isotig08990
<i>suppressor of Hairy wing</i>	2	A	1790-1907	<i>su(Hw)</i>	isotig01368, isotig01367
<i>Suppressor of variegation 3-3</i>	2	A, S	335-422	<i>Su(var)3-3</i>	isotig19729, GFCP6CO01BGKBV
<i>Syntaxin 1A</i>	1	A	4416	<i>Syx1A</i>	isotig04870
<i>TATA box binding protein-related factor 2</i>	2	A	3377-3469	<i>Trf2</i>	isotig01886, isotig01885
<i>TBP-associated factor 1</i>	1	A	5541	<i>Taf1</i>	isotig04746
<i>terribly reduced optic lobes</i>	1	A	690	<i>trol</i>	isotig15574
<i>Trithorax-like</i>	1	A	1267	<i>Trl</i>	isotig11437
<i>warts</i>	1	A	734	<i>wtS</i>	isotig14894
<i>widerborst</i>	1	A	2141	<i>wdb</i>	contig21405

1. Aboïm AN (1945) Développement embryonnainre et post-embryonnaire des gonades normales et agamétiques de *Drosophila melanogaster*. Revue Suisse de Zoologie 3: 53-154.
2. Le Bras S, Van Doren M (2006) Development of the male germline stem cell niche in *Drosophila*. Developmental Biology 294: 92-103.
3. Nelsen OE (1931) Life cycle, sex differentiation, and testis development in *Melanoplus differentialis* (Acrididae, Orthoptera). Journal of Morphology 51: 467-525.

Table S6

Selected developmental process genes identified in the *de novo* *G. bimaculatus* transcriptome.

Process	# Hits	Hit ID (A/S)	Length (range)	Query Gene	Transcriptome Sequence Name(s)
MATERNAL GENES					
ANTERIOR GROUP					
<i>bicoid interacting protein 1</i>	1	A	1040	<i>Bin1</i>	isotig03457
<i>exuperantia</i>	2	A	3152-3225	<i>exu</i>	isotig04765, isotig04764
<i>staufer</i>	3	A	1287-1442	<i>stau</i>	isotig03172, isotig03173, isotig03174
POSTERIOR GROUP					
<i>armitage</i>	1	A	4095	<i>armi</i>	isotig07934
<i>Bruno</i>	1	A	1676	<i>aret</i>	isotig10307
<i>cappuccino</i>	2	A	817-866	<i>capu</i>	isotig06798, isotig06799
<i>fat facets</i>	5	A	1816-3259	<i>faf</i>	isotig01188, isotig01187, isotig01186, isotig01185, isotig01184
<i>Moesin</i>	1	A	4272	<i>Moe</i>	isotig00886
<i>mago nashi</i>	1	A	1021	<i>mago</i>	isotig12375
<i>par-1</i>	1	A	889	<i>par-1</i>	isotig07610
<i>pipsqueak</i>	2	A, S	337-430	<i>psq</i>	isotig19171, GFCP6CO01CETJB
<i>pumilio</i>	3	A, S	412-624	<i>pum</i>	isotig04477, isotig04476, GFJY65E02G1R75
<i>orb</i>	1	A	4765	<i>orb</i>	isotig00462
<i>Rabenosyn-5</i>	1	A	1853	<i>Rbsn-5</i>	isotig09916
<i>staufer</i>	3	A	1287-1442	<i>stau</i>	isotig03172, isotig03173, isotig03174
<i>tudor</i>	2	A	4146-5784	<i>tud</i>	isotig07719, isotig07925
<i>vasa</i>	2	A	765-1146	<i>vas</i>	isotig14543, isotig11874
<i>ypsilon schachtel</i>	1	A	2601	<i>yps</i>	isotig03079
TERMINAL GROUP					
<i>capicua</i>	2	S	314-438	<i>cic</i>	GE8SX9M02IXJOG, GE8SX9M01D8UIJ

<i>corkscrew</i>	1	S	266	<i>csw</i>	GE8SX9M02G96K3
<i>pole hole</i>	1	A	4282	<i>phl</i>	isotig07892
<i>Ras oncogene at 85D</i>	2	A	2078-2467	<i>Ras85D</i>	isotig09494, isotig08979
<i>rolled</i>	1	A	799	<i>rl</i>	isotig14164
<i>torso-like</i>	1	S	174	<i>ts1</i>	GFCP6CO02G92YK
DORSAL GROUP					
<i>cactus</i>	4	A	3168-4301	<i>cact</i>	isotig02364, isotig02362, isotig02363, isotig02361
<i>cappuccino</i>	2	A	817-866	<i>capu</i>	isotig06798, isotig06799
<i>cornichon</i>	1	A	1733	<i>cni</i>	isotig05694
<i>capicua</i>	2	S	314-438	<i>cic</i>	GE8SX9M02IXJOG, GE8SX9M01D8UIJ
<i>dorsal</i>	5	A, S	325-810	<i>dl</i>	isotig14031, GE8SX9M02HRGAV, GFJY65E02GK63W, GFJY65E02FIMPE, GE8SX9M01CGCYQ
<i>Egfr</i>	1	A	1099	<i>Egfr</i>	isotig12088
<i>gastrulation-defective</i>	1	A	862	<i>gd</i>	isotig13529
<i>Myd88</i>	1	A	2079	<i>Myd88</i>	isotig09497
<i>orb</i>	1	A	4765	<i>orb</i>	isotig00462
<i>pelle</i>	2	A	3507-4221	<i>pll</i>	isotig02382, isotig02381
<i>pipe</i>	1	A	6608	<i>pip</i>	isotig07697
<i>spatzle</i>	1	A	2006	<i>spz</i>	isotig09642
<i>squid</i>	1	A	1546	<i>sqd</i>	isotig00544
<i>Toll</i>	1	A	2125	<i>Tl</i>	isotig09438
<i>zucchini</i>	1	A	1455	<i>zuc</i>	isotig00915
ZYGOTICALLY TRANSCRIBED GENES					
<i>cap-n-collar</i>	2	A	1549-2281	<i>cnc</i>	isotig05578, isotig05577
<i>crocodile</i>	2	A	890-966	<i>croc</i>	isotig06650, isotig06649
<i>Tenascin major</i>	7	A, S	200-833	<i>Ten-m</i>	GFJY65E01CUG9F, GE8SX9M01AOG18, GFCP6CO01DGZ87, FQTBZRY01EVWST, GFCP6CO02HATIX, GFCP6CO02G16S1, isotig13797
<i>C-terminal binding protein</i>	3	A, S	239-624	<i>CtBP</i>	isotig16142, GE8SX9M01EF4BJ, FQTBZRY01BYCPR
<i>domeless</i>	1	A	927	<i>dome</i>	isotig12992
<i>eyelid</i>	1	A	2298	<i>osa</i>	isotig09196
<i>ftz transcription factor 1</i>	1	S	397	<i>ftz-f1</i>	GFCP6CO02HU50W
<i>hopscotch</i>	1	A	2719	<i>hop</i>	isotig04276

<i>marelle</i>	1	A	2243	<i>Stat92E</i>	isotig03185
<i>Rpd3</i>	1	A	2212	<i>Rpd3</i>	isotig09325
<i>shuttle craft</i>	1	A	4369	<i>stc</i>	isotig07864
<i>Sir2</i>	1	A	2334	<i>Sir2</i>	contig14671
<i>squid</i>	1	A	1546	<i>sqd</i>	isotig00544

We have amended this subheading title to “Automated annotation using the custom script “Gene Predictor” identifies 14,130 transcriptome sequences as putatively orthologous to D. melanogaster genes.”

"Coding potential of unidentified transcripts" is not very informative. What is the take home message of this section?

*We have amended this subheading title to “Transcripts lacking significant BLAST hits against **nr** may encode functional protein domains.”*

"Analysis of putative orthopteroid-specific sequences" is vague. What type of analysis was conducted? What does the analysis show?

*We have amended this subheading title to “Taxonomic bias of the **nr** database can limit gene identification in de novo assembled transcriptomes.”*

The last paragraph of this section (p. 26) may deserve its own subheading.

We have provided this section with the subheading “Putative orthopteroid-specific sequences contain a high proportion of predicted protein coding domains of unknown function (DUFs).”

Finally, the manuscript a number of grammatical, spelling, and stylistic flaws that should be fixed. Example sentences include: In the Abstract "This database has greatly expanded?" is redundant with previous sentences,

We have removed this redundancy by restructuring the sentence in question.

in the Introduction "An EST project used Sanger sequencing to produce?" is confusingly worded,

The entire section containing this section has been removed in response to this reviewer’s suggestions points #6 and #29.

in the Introduction "Existing genomic resources have thus focused?" the word "thus" brings a different meaning and the "However" in the following sentence also appear to be misleading to the purposes of the sentences, etc.

We have eliminated the words “thus” and “however” from these sentences.

Figure 4 caption refers to "large numbers" but should instead read "numbers in large font".

We have changed this text as suggested.

On page 17, the plural in "showed similarity with these criteria" is confusing. Isn't it just the cutoff at $1e-5$?

Yes, the similarity criteria are those defined in the previous sentence. We have removed this